
Information-theoretic Generalization Analysis for Vector-Quantized VAEs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Encoder-decoder models, which transform input data into latent variables, have
2 achieved significant success in machine learning. While the generalization ability of
3 these models has been theoretically analyzed in supervised learning focusing on the
4 complexity of latent variables, the role of latent variables in generalization and data
5 generation performances are less explored theoretically in unsupervised learning.
6 To address this gap, our study leverages information-theoretic generalization error
7 analysis (IT analysis). Using the supersample setting in recent IT analysis, we
8 demonstrate that the generalization gap for reconstruction loss can be evaluated
9 through mutual information related to the posterior distribution of latent variables
10 conditional on the input data, without relying on the decoder’s information. We
11 also introduce a novel permutation-symmetric supersample setting, which extends
12 the existing IT analysis and shows that regularization of the encoder’s capacity
13 leads to generalization. Finally, we guarantee the Wasserstein distance between the
14 data distribution and the distribution of generated data, offering insights into the
15 model’s data generation capabilities.

16 1 Introduction

17 Encoder-decoder models have achieved significant success in machine learning (Goodfellow et al.,
18 2016). Typically, the encoder extracts information from input data to generate appropriate represen-
19 tations, called latent variables, and the decoder uses these representations to output predictions. In
20 supervised learning, these models are trained by minimizing empirical loss, and regularization of lat-
21 ent variables helps prevent overfitting, improving generalization performance. Many existing studies
22 on encoder-decoder models have focused not only on learned parameters but also on the complexity
23 of latent variables, through principles such as the minimum description length (MDL) (Grnwald
24 et al., 2005), PAC-Bayes (McAllester, 1998), and the information bottleneck (IB) hypothesis (Tishby
25 et al., 2000). More recently, Sefidgaran et al. (2023) theoretically studied latent variable models
26 using the information-theoretic analysis demonstrating that generalization can be characterized by
27 the complexity of the encoder and latent variables without relying on decoder information.

28 Encoder-decoder models are also widely used in unsupervised learning, particularly in deep generative
29 models. When learning these models, we minimize reconstruction loss, which measures the difference
30 between the original data and the regenerated data obtained by compressing data into latent variables
31 by the encoder and regenerating the data by the decoder. Similar to supervised learning, regularization
32 of the latent variables plays a critical role. For example, in variational autoencoder (VAE) (Kingma,
33 2013), in the case of a Gaussian likelihood, the reconstruction loss corresponds to the squared loss,
34 and the regularization term is the Kullback–Leibler (KL) divergence between the prior and posterior
35 distributions of the latent variables. There have been numerous empirical and qualitative studies to
36 explore model performance using the IB hypothesis and rate-distortion theory (Cover & Thomas,

2012) (Alemi et al., 2018; Blau & Michaeli, 2019; Tschannen et al., 2020; Bond-Taylor et al., 2021), but theoretical advances remain limited. Most research has concentrated on encoder and decoder parameters (Epstein & Meir, 2019; Chérif-Abdellatif et al., 2022), leaving a limited understanding of how latent variables contribute to model performance. Although Mbacke et al. (2023) recently introduced PAC-Bayes bounds that use priors and posteriors over the latent variables, their work assumed fixed encoder and decoder parameters, without considering learning these parameters.

Based on the existing research, we provide a theoretical analysis that guarantees unsupervised learning models’ generalization and data generation capabilities, focusing on latent variables. However, simply extending the analysis of VAEs (Mbacke et al., 2023) results in the bounds that depend on learned decoder parameters, obscuring the role of latent variables. Similarly, directly using the information-theoretic analysis from supervised learning (Harutyunyan et al., 2021; Hellström & Durisi, 2022) is challenging due to the difficulty in decoupling the encoder-decoder relationship.

To address these challenges, we propose a novel information-theoretic generalization error bound (Theorem 2) for models with finite latent variables, such as VQ-VAEs (Van Den Oord et al., 2017) (detailed in Sec. 2.1), based on the supersample setting in existing IT analysis incorporating techniques from Mbacke et al. (2023) and Sefidgaran et al. (2023). Furthermore, we introduce a novel permutation-invariant supersample setting, ensuring the generalization gap vanishes as we increase the sample size (Theorem 3). Finally, we provide a guarantee for the data-generating ability by deriving the upper bound on the 2-Wasserstein distance between the data distribution and the distribution of generated data (Theorem 7). These findings provide the first comprehensive theoretical understanding of how encoders and latent variables contribute to generalization and data generation capabilities.

2 Preliminaries

For a random variable (RV) denoted in capital letters, we express its realization with corresponding lowercase letters. Let $p(X)$ denote the distribution of X , and let $p(Y|X)$ represent the conditional distribution of Y given X . We express the expectation of a random variable X as $\mathbb{E}_{p(X)}$ or \mathbb{E}_X . The symbol $I(X; Y)$ represents the mutual information (MI) between X and Y , while $I(X; Y|Z)$ is the conditional MI (CMI) between X and Y given Z . The Kullback–Leibler (KL) divergence between $p(X)$ and $p(Y)$ is denoted $\text{KL}(p(X)||p(Y))$. We further define $[n] = \{1, \dots, n\}$ for $n \in \mathbb{N}$.

2.1 Settings of the latent variable model

This work focuses on encoder-decoder models for unsupervised learning, specifically those with discrete latent spaces, including models such as the vector quantized VAE (VQ-VAE) (Van Den Oord et al., 2017) and its stochastic extensions (Williams et al., 2020; Takida et al., 2022; Sønderby et al., 2017; Roy et al., 2018). Let $\mathcal{X} \subset \mathbb{R}^d$ be the data space and we assume an unknown data generating distribution \mathcal{D} . We express the latent space $\mathcal{Z} \subset \mathbb{R}^{d_z}$, with both \mathcal{X} and \mathcal{Z} equipped with the Euclidean metric $\|\cdot\|$. In the discrete latent space, there are K distinct points, represented as $\mathbf{e} = \{e_j\}_{j=1}^K \in \mathcal{Z}^K$, which are collectively referred to as a codebook learned from the training data, as explained below. Encoder-decoder models consist of two components: an encoder network $f_\phi : \mathcal{X} \rightarrow \mathcal{Z}$ and a decoder network $g_\theta : \mathcal{Z} \rightarrow \mathcal{X}$, parameterized by $\phi \in \Phi \subset \mathbb{R}^{d_e}$ and $\theta \in \Theta \subset \mathbb{R}^{d_d}$, respectively. For a given data point x , the encoder network transforms it into $f_\phi(x)$ and selects the corresponding discrete representation e_j from the codebook \mathbf{e} . The posterior categorical distribution over the index is given as $q(J = j|\mathbf{e}, \phi, x)$ for $j = 1, \dots, K$. We will introduce examples of this distribution later. Using selected latent representation e_J , the decoder network reconstructs the data as $g_\theta(e_J)$. To generate new data, the index J is drawn from a prior distribution, such as a uniform distribution, and the decoder network returns $g_\theta(e_J)$.

Given a training dataset $S = (S_1, \dots, S_n) \in \mathcal{X}^n$, where each data point $S_m \in \mathcal{X}$ is drawn i.i.d from \mathcal{D} , we jointly learn the parameters of the encoder, decoder, and the codebook. We denote the set of parameters as $W := \{\mathbf{e}, \phi, \theta\} \in \mathcal{W} := \mathcal{Z}^K \times \Phi \times \Theta$. We assume that these parameters are learned using a randomized algorithm and the learning process is represented by the conditional distribution $\mathbf{e}, \phi, \theta \sim q(\mathbf{e}, \phi, \theta|S)$. The learning algorithm typically minimizes the reconstruction loss. For a given data point x and the corresponding latent variable e_j , the quality of the reconstructed data is measured by a loss function $l(x, g_\theta(e_j))$, where $l : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+$. Then the reconstruction loss for input x and parameter w is defined as $l_0 : \mathcal{W} \times \mathcal{X} \rightarrow \mathbb{R}$, $l_0(w, x) := \mathbb{E}_{q(J|\mathbf{e}, \phi, x)} l(x, g_\theta(e_J))$. In this work, we focus on the squared distance for the loss function l , so we aim to minimize $l_0(w, x) := \mathbb{E}_{q(J|\mathbf{e}, \phi, x)} \|x - g_\theta(e_J)\|^2$ over the training dataset $x \in S$.

91 Finally, we provide examples of $q(J|\mathbf{e}, \phi, x)$. The original VQ-VAE (Van Den Oord et al., 2017)
 92 used the deterministic process

$$q(J = j|\mathbf{e}, \phi, x) = \begin{cases} 1 & \text{for } j = \arg \min_{k \in [K]} \|f_\phi(x) - e_k\|, \\ 0 & \text{otherwise,} \end{cases}$$

93 using the distance between the outputs of the encoder and the codebook. Recently, stochastic selection
 94 methods have gained popularity. For instance, Williams et al. (2020) proposed the distribution

$$q(J = j|\mathbf{e}, \phi, x) \propto \exp(-\beta \|f_\phi(x) - e_j\|^2), \quad (1)$$

95 where the softmax function is used, and $\beta \in \mathbb{R}^+$ is a temperature parameter that controls the level of
 96 stochasticity. Beyond this, using stochastic encoders has become common in several other works,
 97 including Sønderby et al. (2017); Roy et al. (2018); Takida et al. (2022).

98 2.2 Information-theoretic generalization error analysis

99 We now briefly outline the IT analysis using the supersample that we utilize in our study (Steinke
 100 & Zakynthinou, 2020; Harutyunyan et al., 2021; Hellström & Durisi, 2022). Note that the existing
 101 IT analysis is used for supervised learning, the notation of this section is slightly different from
 102 our main results in Sec.3. Let \mathcal{X} be the domain of data and suppose \mathcal{D} represents an *unknown* data
 103 distribution. Consider $\tilde{X} \in \mathcal{X}^{n \times 2}$ as an $n \times 2$ matrix, where each entry is drawn i.i.d. from \mathcal{D} .
 104 We refer to this matrix as a *supersample*. Each column of \tilde{X} has indexes $\{0, 1\}$ associated with
 105 $U = (U_1, \dots, U_n) \sim \text{Uniform}(\{0, 1\}^n)$ independent of \tilde{X} . We denote the m -th row as \tilde{X}_m with
 106 entries $(\tilde{X}_{m,0}, \tilde{X}_{m,1})$. In this setting, we consider $\tilde{X}_U := (\tilde{X}_{m,U_m})_{m=1}^n$ as the training dataset and
 107 $\tilde{X}_{\bar{U}} := (\tilde{X}_{m,\bar{U}_m})_{m=1}^n$ as the test dataset, where $\bar{U}_m = 1 - U_m$. We consider a randomized algorithm
 108 $\mathcal{A} : \mathcal{X}^n \rightarrow \mathcal{W}$, where $w \in \mathcal{W} \subset \mathbb{R}^{d_w}$ is a parameter. Given a training dataset S , the learning
 109 algorithm can be characterized by $q(W|S)$. We evaluate the quality of the learning algorithm using
 110 the loss function $l : \mathcal{W} \times \mathcal{X} \rightarrow [0, 1]$, where $l(\mathcal{A}(s), x)$ for fixed $S = s$ and $X = x$. With these
 111 notations, $l(\mathcal{A}(\tilde{X}_U), \tilde{X})$ denotes the $n \times 2$ loss matrix obtained by applying $l(\mathcal{A}(\tilde{X}_U), \cdot)$ elementwise
 112 to \tilde{X} . In this setting, we can see that $\hat{L}_{\tilde{X}} := \frac{1}{n} \sum_{m=1}^n l(\mathcal{A}(\tilde{X}_U), \tilde{X}_{m,U_m})$ corresponds to the training
 113 error and $L_{\tilde{X}} := \frac{1}{n} \sum_{m=1}^n l(\mathcal{A}(\tilde{X}_U), \tilde{X}_{m,\bar{U}_m})$ corresponds to the test error. The described settings
 114 called the **supersample setting** lead to the following generalization error bound:

115 **Theorem 1** (Hellström & Durisi (2022)). *Under the supersample setting, we have*

$$|\mathbb{E}_{\tilde{X}, U}(L_{\tilde{X}} - \hat{L}_{\tilde{X}})| \leq \sqrt{\frac{2}{n} I(l(\mathcal{A}(\tilde{X}_U), \tilde{X}); U|\tilde{X})}.$$

116 3 Generalization of the reconstruction loss

117 This section aims to analyze the generalization capability of encoder-decoder models using IT
 118 analysis. We define the generalization error of the reconstruction loss as follows:

$$\text{gen}(n, \mathcal{D}) := \left| \mathbb{E}_{S, \tilde{X}} \mathbb{E}_{q(\mathbf{e}, \phi, \theta|S)} \left(\mathbb{E}_{q(J|\mathbf{e}, \phi, X)} l(X, g_\theta(e_J)) - \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m|\mathbf{e}, \phi, S_m)} l(S_m, g_\theta(e_{J_m})) \right) \right|.$$

119 To proceed with the analysis, we assume the following condition regarding the data space:

120 **Assumption 1.** *There exists a positive constant Δ such that $\sup_{x, x' \in \mathcal{X}} \|x - x'\| < \Delta^{1/2}$.*

121 This assumption implies that for any x and e_j and θ , the loss function $l(x, g_\theta(e_j))$ is bounded by Δ .

122 We now restate the settings from Sec. 2.1 under the supersample framework. Given a super-
 123 sample $\tilde{X} := (\tilde{X}_0, \tilde{X}_1) \in \mathcal{X}^{n \times 2}$, define $\tilde{X}_U := (\tilde{X}_{m,U_m})_{m=1}^n$ as the training dataset and
 124 $\tilde{X}_{\bar{U}} := (\tilde{X}_{m,\bar{U}_m})_{m=1}^n$ as the test dataset. Then treating $l_0(w, x) := \mathbb{E}_{q(J|\mathbf{e}, \phi, x)} \|x - g_\theta(e_J)\|^2$
 125 as l in Sec 2.2, we can directly apply the generalization bound in Theorem 1. We refer to this general-
 126 ization bound as the **naive IT-bound** (See Appendix B for the formal statement.). As discussed in
 127 Appendix B, the naive IT-bound does not clearly capture the role of the learned representation e_J
 128 in generalization because the CMI term is entangled with both the learning of $W = \{\mathbf{e}, \phi, \theta\}$ and
 129 posterior distribution $q(J|\mathbf{e}, \phi, x)$. This section aims to extend the naive IT analysis to the bound that
 130 captures the role of representation.

131 **3.1 The generalization error under the existing supersample setting**

132 We introduce the notations of the joint distributions. Given a super sample \tilde{X} , we define
 133 $q(\tilde{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}_{\bar{U}}) = \prod_{m=1}^n q(\tilde{J}_m|\mathbf{e}, \phi, \tilde{X}_{m, \bar{U}_m})$ and $q(\mathbf{J}|\mathbf{e}, \phi, \tilde{X}_U) = \prod_{m=1}^n q(J_m|\mathbf{e}, \phi, \tilde{X}_{m, U_m})$ and
 134 $q(\tilde{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}) = q(\tilde{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \phi, \tilde{X}_{\bar{U}}, \tilde{X}_U) = q(\tilde{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}_{\bar{U}})q(\mathbf{J}|\mathbf{e}, \phi, \tilde{X}_U)$.

135 Following is our first main result, the proof is shown in Appendix C.

136 **Theorem 2.** *Under Assumption 1 and the supersample setting, we have*

$$\text{gen}(n, \mathcal{D}) \leq 2\Delta \sqrt{\frac{I(\tilde{\mathbf{J}}; U|\mathbf{e}, \phi, \tilde{X}) + \mathbb{E}_{\tilde{X}, U} \mathbb{E}_{q(\mathbf{e}, \phi, \theta|\tilde{X}_U)} \text{KL}(\mathbf{Q}|\mathbf{P})}{n}} + \frac{\Delta}{\sqrt{n}},$$

137 where the CMI is defined as

$$I(\tilde{\mathbf{J}}; U|\mathbf{e}, \phi, \tilde{X}) = \mathbb{E}_{\tilde{X}, U} \mathbb{E}_{q(\mathbf{e}, \phi|\tilde{X}_U)} \text{KL}(q(\tilde{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}) \parallel \mathbb{E}_{U'} q(\tilde{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \phi, \tilde{X}_{\bar{U}'}, \tilde{X}_{U'})).$$

138 The distributions of KL divergence are defined as

$$\mathbf{Q} := q(\mathbf{e}, \phi, \theta|\tilde{X}_U) \prod_{m=1}^n q(J_m|\mathbf{e}, \phi, \tilde{X}_{m, U_m}), \quad \mathbf{P} := q(\mathbf{e}, \phi, \theta|S) \prod_{m=1}^n q(J_m|\mathbf{e}, \phi),$$

139 and $q(J_m|\mathbf{e}, \phi)$ is any prior distribution that does not depend on the training data.

140 The bound does not depend on the decoder’s information; This means that even if a complex decoder
 141 network is used to reduce reconstruction loss, it does not worsen the generalization gap. The CMI
 142 and KL terms are influenced solely by the posterior distribution of the latent variables, conditioned
 143 on the learned ϕ and \mathbf{e} .

144 **The role of the representation in our bound:** Denoting $\tilde{X}_U = S = (S_1, \dots, S_n)$, the KL
 145 divergence term can be rewritten as

$$\frac{\mathbb{E}_{q(\mathbf{e}, \phi, \theta|S)} \text{KL}(\mathbf{Q}|\mathbf{P})}{n} = \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi|S)} \text{KL}(q(J_m|\mathbf{e}, \phi, S_m) \parallel q(J_m|\mathbf{e}, \phi)).$$

146 This is referred to as **the empirical KL divergence** in Mbacke et al. (2023), which is often used as
 147 the regularization in the variational inference. For the CMI term, since $\tilde{X}_{\bar{U}}$ are n i.i.d RVs from \mathcal{D} ,
 148 we express each $\tilde{X}_{\bar{U}}$ as X and then, we have the following relation, see Appendix D.1 for its proof;

$$I(\tilde{\mathbf{J}}; U|\mathbf{e}, \phi, \tilde{X}) \leq nI(e_J; X|\mathbf{e}, \phi) + \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi|S)} \frac{1}{n} \sum_{m=1}^n \text{KL}(q(J_m|\mathbf{e}, \phi, S_m) \parallel q(J_m|\mathbf{e}, \phi)). \quad (2)$$

149 This upper bound is characterized by the mutual information (MI), which is commonly used in the IB
 150 hypothesis, and the empirical KL divergence term. These are popular empirical evaluation metrics.
 151 However, as discussed in Sefidgaran et al. (2023), these terms do not vanish as $n \rightarrow \infty$. Therefore,
 152 the following discussion suggests that utilizing the symmetry of the prior distribution (concerning the
 153 supersample) is important to address such issues.

154 **The dependency on the sample size:** Next, we study the dependency of the CMI and KL term on n
 155 in Theorem 2. The CMI term is similar to the fCMI term from existing IT analysis (Harutyunyan et al.,
 156 2021), but here, the conditioning on all other parameters distinguishes it from typical fCMI bounds,
 157 see Sec. D.3 for the detailed discussion. Since the latent space is discrete, we have $I(\tilde{\mathbf{J}}; U|\mathbf{e}, \phi, \tilde{X}) \leq$
 158 $2n \log K$, ensuring that the bound is always finite, though it may be vacuous. When using the
 159 deterministic decoder $f_\phi : \mathcal{X} \rightarrow [K]$, we can directly use Theorem 8 in Hellström & Durisi (2022);
 160 if f_ϕ belongs to a class of functions that has a finite Natarajan-dimension, then $I(\tilde{\mathbf{J}}; U|\mathbf{e}, \phi, \tilde{X}) =$
 161 $\mathcal{O}(\log n)$, see Appendix D.2 for the details. Thus, by regularizing the encoder model’s capacity, the
 162 first term inside the square root in Theorem 2 scales as $\mathcal{O}(\log n/n)$. Comparing this with Eq. (2),
 163 where $I(e_J; X|\phi)$ does not vanish as $n \rightarrow \infty$, this highlights the importance of using symmetry in
 164 the prior distribution for supersamples to achieve meaningful bounds, as discussed in Sefidgaran et al.
 165 (2023). For a stochastic encoder, like in Eq.(1), regularizing the encoder network’s capacity similarly
 166 bounds the CMI (see Appendix F and Theorem 4 in the below).

167 Regarding the empirical KL term, it is larger than the CMI term as seen in Eq. (2), and it does not
 168 necessarily vanish as $n \rightarrow \infty$, as pointed out in Geiger & Koch (2019) and Sefidgaran et al. (2023).
 169 As discussed in Appendix C.3, this arises from the limited flexibility of the supersample setting,
 170 which motivated the introduction of the novel supersample setting in Sec. 3.2.

171 **3.2 Generalization under the permutation symmetry settings**

172 As discussed in Sec. 3.1, the existing supersample setting leads to an empirical KL term that does not
 173 necessarily vanish as n increases. As discussed in Sefidgaran et al. (2023), the existing supersample
 174 setting utilizes the specific symmetry of the test and training dataset (they referred to it as type-1
 175 symmetry) and demonstrated that such symmetry is insufficient to analyze latent variable models.
 176 We extend their results by introducing a new symmetry, which eliminates the empirical KL term.

177 To establish this new symmetry, let us denote a random permutation of $[2n]$ as $\mathbf{T} = \{T_1, \dots, T_{2n}\}$,
 178 where each permutation appears with uniform probability, $P(\mathbf{T}) = 1/(2n)!$. Given a supersample
 179 $\tilde{X} = (\tilde{X}_1, \dots, \tilde{X}_{2n}) \in \mathcal{X}^{2n}$, a set of $2n$ random variables drawn i.i.d from \mathcal{D} , we reorder the samples
 180 using \mathbf{T} expressed as $\tilde{X}_{\mathbf{T}} = (\tilde{X}_{T_1}, \dots, \tilde{X}_{T_{2n}})$. The first n samples $(\tilde{X}_{T_1}, \dots, \tilde{X}_{T_n})$ are used for
 181 the test dataset and the remaining n samples $(\tilde{X}_{T_{n+1}}, \dots, \tilde{X}_{T_{2n}})$ are used for the training dataset.
 182 We further express $\mathbf{T} = \{\mathbf{T}_0, \mathbf{T}_1\}$ and $\tilde{X}_{\mathbf{T}_0} = (\tilde{X}_{T_1}, \dots, \tilde{X}_{T_n})$ and $\tilde{X}_{\mathbf{T}_1} = (\tilde{X}_{T_{n+1}}, \dots, \tilde{X}_{T_{2n}})$
 183 represent the test and training dataset respectively. Unlike the existing supersample setting discussed
 184 in Sec. 2.2, where U_m are independent, the components of \mathbf{T} are dependent.

185 We express the joint distribution as follows; $q(\bar{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}_0}) = \prod_{m=1}^n q(\bar{J}_m|\mathbf{e}, \phi, \tilde{X}_{T_m})$,
 186 $q(\mathbf{J}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}_1}) = \prod_{m=1}^n q(J_m|\mathbf{e}, \phi, \tilde{X}_{T_{n+m}})$, and $q(\tilde{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}) = q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}_0}, \tilde{X}_{\mathbf{T}_1}) =$
 187 $q(\bar{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}_0})q(\mathbf{J}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}_1})$. We refer to these notations and assumptions as **the permutation**
 188 **symmetric (supersample) setting**. Following is our main result, the proof is shown in Appendix E;

189 **Theorem 3.** *Under Assumptions 1 and the permutation symmetric setting, we have*

$$\text{gen}(n, \mathcal{D}) \leq 4\Delta \mathbb{E}_X \sqrt{\frac{I(\tilde{\mathbf{J}}; \mathbf{T}|\mathbf{e}, \phi, \tilde{X})}{n}} + \frac{2\Delta}{\sqrt{n}},$$

190 where the CMI is defined as

$$I(\tilde{\mathbf{J}}; \mathbf{T}|\mathbf{e}, \phi, \tilde{X}) = \mathbb{E}_{\tilde{X}, \mathbf{T}} \mathbb{E}_{q(\mathbf{e}, \phi | \tilde{X}_{\mathbf{T}_1})} \text{KL}(q(\tilde{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}) \| \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}_0}, \tilde{X}_{\mathbf{T}_1})). \quad (3)$$

191 As shown, the empirical KL term is eliminated, and a new CMI term, Eq. (3), emerges, which
 192 leverages the symmetry of index \mathbf{T} in the prior distribution. We will show that this CMI term will
 193 vanish as $n \rightarrow \infty$, thus Theorem 3 successfully characterizes the generalization. As discussed in
 194 Sec. 3.1, when using the sufficiently regularized deterministic decoder $f_\phi : \mathcal{X} \rightarrow [K]$, this CMI
 195 scales as $\mathcal{O}(\log n)$, and thus, the bound behaves as $\mathcal{O}(\sqrt{\log n/n})$. See AppendixD.2 for more details.

196 To analyze the role of the capacity of stochastic encoders like Eq. (1), we extend Theorem 3 by
 197 incorporating the concept of *metric entropy*. Assume $q(J|\mathbf{e}, \phi, x) = q(J|\mathbf{e}, f_\phi(x))$. Conditioned on ϕ ,
 198 let \mathcal{F} be the encoder function class equipped with the metric $\|\cdot\|_\infty$. Given $x^n := (x_1, \dots, x_n) \in \mathcal{X}^n$,
 199 define the pseudo-metric d_n on \mathcal{F} as $d_n(f, g) := \max_{i \in [n]} \|f(x_i) - g(x_i)\|_\infty$ for $f, g \in \mathcal{F}$. The
 200 δ -covering number of \mathcal{F} with respect to d_n is denoted as $\mathcal{N}(\delta, \mathcal{F}, x^n)$, and we define $\mathcal{N}(\delta, \mathcal{F}, n) :=$
 201 $\sup_{x^n \in \mathcal{X}^n} \mathcal{N}(\delta, \mathcal{F}, x^n)$.

202 **Theorem 4.** *Assume that there exists a positive constant Δ_z such that $\sup_{z, z' \in \mathcal{Z}} \|z - z'\| < \Delta_z$.*
 203 *Then, when using Eq. (1) and under the same setting as Theorem 3, for any $\delta \in (0, 1]$, we have*

$$\text{gen}(n, \mathcal{D}) \leq \Delta \sqrt{8\beta n \delta \Delta_z} + 4\Delta \sqrt{\frac{\log \mathcal{N}(\delta, \mathcal{F}, 2n)}{n}} + \frac{2\Delta}{\sqrt{n}}.$$

204 In the proof, we first approximate $f_\phi(x)$ using δ -cover of \mathcal{F} , leading to an approximation error in the
 205 first term. Then the CMI of the δ -cover is bounded by the metric entropy. See Appendix F for the
 206 complete proof, including a more general stochastic encoder beyond Eq. (1). When \mathcal{F} is sufficiently
 207 regularized (such as with Natarajan dimension with margin, see AppendixF for details), the metric
 208 entropy scales as $\mathcal{O}(\log(n/\delta))$, and by setting $\delta = \mathcal{O}(1/n^2)$, we achieve $\text{gen}(n, \mathcal{D}) = \mathcal{O}(\sqrt{\log n/n})$.
 209 This result demonstrates that regularizing the encoder’s capacity leads to better generalization.

210 We are often interested in the data generation capabilities rather than generalization under recon-
 211 struction loss. Specifically, the goal is to generate realistic data by sampling from the latent variable
 212 distribution and passing it through the decoder. We aim for the distribution of generated data to closely

213 approximate the true data distribution. In Theorem 7 of Appendix G, we provide an upper bound
214 on the Wasserstein distance between the true data distribution and the generated data distribution
215 obtained from the pushforward of the prior distribution over latent variables.

216 Theorems 2, 3 (and 7) offer important insights into the roles of the encoder and decoder. To reduce
217 the reconstruction loss on test data and improve data generation capabilities, it is desirable to use a
218 complex decoder, as it can lower the reconstruction loss without increasing the KL or CMI terms,
219 regardless of the sample size. However, using the complex encoder increases the KL and CMI,
220 requiring careful adjustment according to the sample size. This characteristic is specific to latent
221 variable models, highlighting the critical role of the latent variables as the regularization.

222 3.3 Comparison with existing bounds

223 Here we compare our bounds with existing work. Theorem 2 resembles the results of Mbacke et al.
224 (2023) since both bounds include the empirical KL term in the upper bounds, and the posterior
225 distribution corresponds to the variational posterior distribution. The key difference is that Mbacke
226 et al. (2023) assumed fixed encoder and decoder parameters, whereas our analysis incorporates the
227 learning process under the assumption of finite latent space and squared reconstruction loss. A further
228 distinction is that their generalization bound does not go to 0 as $n \rightarrow \infty$ due to two reasons; the
229 presence of the empirical KL term, which we address in Theorem 3 using permutation symmetry.
230 Our technique can be regarded as developing the appropriate prior distribution in PAC-Bayes bound.
231 The second reason is the presence of the average distance $\frac{1}{n} \sum_{m=1}^n \mathbb{E}_X \|X - S_m\|$, which is inherent
232 to the data distribution and may not vanish as $n \rightarrow \infty$. Our use of the squared loss in the analysis
233 mitigates this problematic term, as detailed in Appendix, C. Our proof techniques are motivated from
234 Sefidgaran et al. (2023). However, we could not directly apply their methods, as the reconstruction
235 loss reuses input data, unlike in classification settings. We resolve this by combining the data
236 regeneration technique in the proof of Mbacke et al. (2023). Additionally, we introduced a new
237 permutation symmetric setting, leading to a bound that controls mutual information in Theorem 3.

238 Existing analyses based related to the IB hypothesis (Vera et al., 2018; Hafez-Kolahi et al., 2020;
239 Kawaguchi et al., 2023; Vera et al., 2023) assume both the latent variables and data are discrete,
240 and their bounds explicitly depend on the latent space size or show exponential dependency on the
241 MI. In contrast, we only assume discrete latent variables and the resulting bound does not explicitly
242 depend on the number of discrete states nor exhibit exponential dependency on MI. We believe that
243 our technique can be extended to continuous latent variables, which we leave for future research.

244 4 Conclusion and limitations

245 We provided the first comprehensive analysis of the generalization and data generation capabilities of
246 encoder-decoder models in unsupervised learning based on the IT analysis. Our work highlights the
247 role of encoder capacity and the posterior distribution of latent variables through the use of a novel
248 permutation-symmetric supersample setting. However, our analysis has two key limitations. First, it
249 assumes a discrete latent space, limiting its applicability to models like VAEs with continuous latent
250 variables. Second, it relies on the squared loss for reconstruction. Addressing these limitations in
251 future work will be crucial for developing a more accurate understanding of encoder-decoder models.

252 References

- 253 Alexander Alemi, Ben Poole, Ian Fischer, Joshua Dillon, Rif A Saurous, and Kevin Murphy. Fixing a
254 broken elbow. In *International conference on machine learning*, pp. 159–168. PMLR, 2018.
- 255 Noga Alon, Shai Ben-David, Nicolo Cesa-Bianchi, and David Haussler. Scale-sensitive dimensions,
256 uniform convergence, and learnability. *Journal of the ACM (JACM)*, 44(4):615–631, 1997.
- 257 Peter L Bartlett and Wolfgang Maass. Vapnik-chervonenkis dimension of neural nets. *The handbook
258 of brain theory and neural networks*, pp. 1188–1192, 2003.
- 259 S. Bendavid, N. Cesabianchi, D. Haussler, and P.M. Long. Characterizations of learnability for classes
260 of $[0, \dots, n)$ -valued functions. *Journal of Computer and System Sciences*, 50(1):74–86, 1995. ISSN
261 0022-0000. doi: <https://doi.org/10.1006/jcss.1995.1008>. URL <https://www.sciencedirect.com/science/article/pii/S0022000085710082>.

- 263 Yochai Blau and Tomer Michaeli. Rethinking lossy compression: The rate-distortion-perception
264 tradeoff. In *International Conference on Machine Learning*, pp. 675–685. PMLR, 2019.
- 265 Sam Bond-Taylor, Adam Leach, Yang Long, and Chris G Willcocks. Deep generative modelling:
266 A comparative review of vaes, gans, normalizing flows, energy-based and autoregressive models.
267 *IEEE transactions on pattern analysis and machine intelligence*, 44(11):7327–7347, 2021.
- 268 Badr-Eddine Chérif-Abdellatif, Yuyang Shi, Arnaud Doucet, and Benjamin Guedj. On pac-bayesian
269 reconstruction guarantees for vaes. In *International conference on artificial intelligence and
270 statistics*, pp. 3066–3079. PMLR, 2022.
- 271 T. M. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley & Sons, 2012.
- 272 Amit Daniely, Sivan Sabato, Shai Ben-David, and Shai Shalev-Shwartz. Multiclass learnability and
273 the erm principle. In *Proceedings of the 24th Annual Conference on Learning Theory*, pp. 207–232.
274 JMLR Workshop and Conference Proceedings, 2011.
- 275 Devdatt P Dubhashi and Desh Ranjan. Balls and bins: A study in negative dependence. *BRICS
276 Report Series*, 3(25), 1996.
- 277 Baruch Epstein and Ron Meir. Generalization bounds for unsupervised and semi-supervised learning
278 with autoencoders. *arXiv preprint arXiv:1902.01449*, 2019.
- 279 Bernhard C. Geiger and Tobias Koch. On the information dimension of stochastic processes. *IEEE
280 Transactions on Information Theory*, 65(10):6496–6518, 2019. doi: 10.1109/TIT.2019.2922186.
- 281 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. [http:
282 //www.deeplearningbook.org](http://www.deeplearningbook.org).
- 283 Lee-Ad Gottlieb, Aryeh Kontorovich, and Robert Krauthgamer. Efficient classification for metric
284 data. *IEEE Transactions on Information Theory*, 60(9):5750–5759, 2014.
- 285 Peter D. Grnwald, In Jae Myung, and Mark A. Pitt. *Advances in Minimum Description Length: Theory
286 and Applications (Neural Information Processing)*. The MIT Press, 2005. ISBN 0262072629.
- 287 Yann Guermeur. Vc theory of large margin multi-category classifiers. *Journal of Machine Learning
288 Research*, 8(85):2551–2594, 2007.
- 289 Yann Guermeur. Lp-norm sauer–shelah lemma for margin multi-category classifiers. *Journal of
290 Computer and System Sciences*, 89:450–473, 2017. ISSN 0022-0000.
- 291 Yann Guermeur. Combinatorial and structural results for gamma-psi-dimensions. *arXiv preprint
292 arXiv:1809.07310*, 2018.
- 293 Hassan Hafez-Kolahi, Shohreh Kasaei, and Mahdiyeh Soleymani-Baghshah. Sample complexity of
294 classification with compressed input. *Neurocomputing*, 415:286–294, 2020. ISSN 0925-2312.
- 295 H. Harutyunyan, M. Raginsky, G. Ver Steeg, and A. Galstyan. Information-theoretic generalization
296 bounds for black-box learning algorithms. In *Advances in Neural Information Processing Systems*,
297 pp. 24670–24682, 2021.
- 298 F. Hellström and G. Durisi. A new family of generalization bounds using samplewise evaluated CMI.
299 In *Advances in Neural Information Processing Systems*, 2022.
- 300 Ying Jin. Upper bounds on the natarajan dimensions of some function classes. In *2023 IEEE
301 International Symposium on Information Theory (ISIT)*, pp. 1020–1025. IEEE, 2023.
- 302 Kumar Joag-Dev and Frank Proschan. Negative association of random variables with applications.
303 *The Annals of Statistics*, pp. 286–295, 1983.
- 304 Kenji Kawaguchi, Zhun Deng, Xu Ji, and Jiaoyang Huang. How does information bottleneck help
305 deep learning? In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan
306 Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine
307 Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 16049–16096. PMLR,
308 23–29 Jul 2023.

- 309 Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 310 Sokhna Diarra Mbacke, Florence Clerc, and Pascal Germain. Statistical guarantees for variational
311 autoencoders using pac-bayesian theory. *Advances in Neural Information Processing Systems*, 36,
312 2023.
- 313 David A McAllester. Some pac-bayesian theorems. In *Proceedings of the eleventh annual conference*
314 *on Computational learning theory*, pp. 230–234, 1998.
- 315 Aurko Roy, Ashish Vaswani, Arvind Neelakantan, and Niki Parmar. Theory and experiments on
316 vector quantized autoencoders. *arXiv preprint arXiv:1805.11063*, 2018.
- 317 Milad Sefidgaran, Abdellatif Zaidi, and Piotr Krasnowski. Minimum description length and gen-
318 eralization guarantees for representation learning. *Advances in Neural Information Processing*
319 *Systems*, 36, 2023.
- 320 Casper Kaae Sønderby, Ben Poole, and Andriy Mnih. Continuous relaxation training of discrete
321 latent variable image models. In *Bayesian DeepLearning workshop, NIPS*, volume 201, 2017.
- 322 T. Steinke and L. Zakynthinou. Reasoning About Generalization via Conditional Mutual Information.
323 In *Proceedings of Thirty Third Conference on Learning Theory*, volume 125, pp. 3437–3452, 2020.
- 324 Yuhta Takida, Takashi Shibuya, Weihsiang Liao, Chieh-Hsin Lai, Junki Ohmura, Toshimitsu Uesaka,
325 Naoki Murata, Shusuke Takahashi, Toshiyuki Kumakura, and Yuki Mitsufuji. SQ-VAE: Variational
326 Bayes on discrete representation with self-annealed stochastic quantization. In Kamalika Chaudhuri,
327 Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of*
328 *the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine*
329 *Learning Research*, pp. 20987–21012. PMLR, 17–23 Jul 2022.
- 330 Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv*
331 *preprint physics/0004057*, 2000.
- 332 Michael Tschannen, Josip Djolonga, Paul K. Rubenstein, Sylvain Gelly, and Mario Lucic. On mutual
333 information maximization for representation learning. In *International Conference on Learning*
334 *Representations*, 2020.
- 335 Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in*
336 *neural information processing systems*, 30, 2017.
- 337 Matias Vera, Pablo Piantanida, and Leonardo Rey Vega. The role of the information bottleneck in
338 representation learning. In *2018 IEEE International Symposium on Information Theory (ISIT)*, pp.
339 1580–1584, 2018. doi: 10.1109/ISIT.2018.8437679.
- 340 Matias Vera, Leonardo Rey Vega, and Pablo Piantanida. The role of mutual information in variational
341 classifiers. *Machine Learning*, 112(9):3105–3150, 2023.
- 342 Will Williams, Sam Ringer, Tom Ash, David MacLeod, Jamie Dougherty, and John Hughes. Hi-
343 erarchical quantized autoencoders. *Advances in Neural Information Processing Systems*, 33:
344 4524–4535, 2020.

345 A Auxiliary definitions and lemmas

346 Here we define the Wasserstein distance. Given a metric $d(\cdot, \cdot)$ and probability distributions p and q
347 on \mathcal{X} , let $\Pi(p, q)$ denote the set of all couplings of p and q . The 2-Wasserstein distance is defined as:

$$W_2(p, q) = \sqrt{\inf_{\rho \in \Pi} \int_{\mathcal{X} \times \mathcal{X}} d(x, x')^2 d\rho(x, x')}.$$

348 In this work, we use the Euclidean metric $|\cdot|$ as $d(\cdot, \cdot)$.

349 We also rely on the following type of exponential moment inequality, which is often used in the proof
 350 of McDiarmid's inequality. A function $f : \mathcal{X}^n \rightarrow \mathbb{R}$ has the bounded differences property if for some
 351 nonnegative constants c_1, \dots, c_n , the following holds for all i :

$$\sup_{x_1, \dots, x_n, x'_i \in \mathcal{X}} |f(x_1, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c_i, \quad 1 \leq i \leq n.$$

352 Assuming X_1, \dots, X_n are independent random variables taking values in \mathcal{X} , we have the following
 353 lemma:

354 **Lemma 1** (Used in the proof of McDiarmid's inequality). *Given a function f with the bounded*
 355 *differences property, for any $t \in \mathbb{R}$, we have:*

$$\mathbb{E} \left[e^{t(f(X_1, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X_n)])} \right] \leq e^{\frac{t^2}{8} \sum_{i=1}^n c_i^2}.$$

356 B Discussion about the Naive IT bound

357 As discussed in Sec 3, by applying the existing IT analysis bound in Theorem 1, we can derive a
 358 naive IT bound for the reconstruction loss as follows:

359 **Theorem 5.** *Under Assumption 1 and the supersample setting, we have*

$$\text{gen}(n, \mathcal{D}) \leq \Delta \sqrt{\frac{2}{n} I(l_0(W, \tilde{X}); U | \tilde{X})}.$$

360 where $l_0(w, x) := \mathbb{E}_{q(\mathbf{J} | \mathbf{e}, \phi, x)} \|x - g_\theta(e_{\mathbf{J}})\|^2$ and $W = \{\mathbf{e}, \phi, \theta\} \sim q(\mathbf{e}, \phi, \theta | \tilde{X}_U)$.

361 *Proof.* Given that the loss is bounded by $[0, \Delta]$, it follows a Δ -subGaussian property. Thus, using
 362 Theorem 1, we obtain the result. \square

363 It is important to note that this upper bound is characterized by the CMI $I(l_0(W, \tilde{X}); U | \tilde{X})$. This
 364 CMI depends on the decoder and encoder information, distinguishing it from the results presented in
 365 our main Theorems 2 and 3, which do not require the decoder's information.

366 To clarify this distinction, let us introduce the necessary notation. Following the notation in Sec. 3.1,
 367 we define the regenerated data as:

$$\tilde{Y} := (g_\theta(e_{\tilde{\mathbf{J}}_1}), \dots, g_\theta(e_{\tilde{\mathbf{J}}_n}), g_\theta(e_{\mathbf{J}_1}), \dots, g_\theta(e_{\mathbf{J}_n})) = g_\theta(e_{\tilde{\mathbf{J}}}),$$

368 which represents the elementwise application of the decoder $g_\theta(e_{(\cdot)})$ to the selected index $\tilde{\mathbf{J}}$ on \tilde{X}
 369 (Recall that $q(\tilde{\mathbf{J}} | \mathbf{e}, \phi, \tilde{X}) = q(\tilde{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, \tilde{X}_{\tilde{V}}, \tilde{X}_U) = q(\tilde{\mathbf{J}} | \mathbf{e}, \phi, \tilde{X}_{\tilde{V}}) q(\mathbf{J} | \mathbf{e}, \phi, \tilde{X}_U)$).

370 Under these notations, we have the following relations:

$$I(l_0(W, \tilde{X}); U | \tilde{X}) \leq I(\tilde{Y}; U | \tilde{X}) \leq I(\theta; U | \tilde{X}) + I(e_{\tilde{\mathbf{J}}}; U | \tilde{X}, \theta)$$

371 where the first inequality is obtained by the data processing inequality (DPI) and the second inequality
 372 is obtained by the chain rule of CMI and the DPI. This result demonstrates that the decoder information
 373 cannot be eliminated from the naive IT bound, which clarifies the fundamental difference compared
 374 to our result (Theorems 2 and 3).

375 C Proof of Theorem 2

376 We express $q(\tilde{\mathbf{J}} | \mathbf{e}, \phi, \tilde{X}) = q(\tilde{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, \tilde{X}_{\tilde{V}}, \tilde{X}_U) = q(\tilde{\mathbf{J}} | \mathbf{e}, \phi, \tilde{X}_{\tilde{V}}) q(\mathbf{J} | \mathbf{e}, \phi, \tilde{X}_U)$. Hereinafter, we
 377 simplify the notation by expressing \tilde{X} as X . For simplification in the proof, we omit the absolute value
 378 operation. The reverse bound can be proven in a similar manner. We first express the generalization

379 error of the reconstruction loss using the supersample as follows

$$\begin{aligned}
& \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\bar{J}_m | \mathbf{e}, \phi, X_m, \bar{U}_m)} q(\mathbf{e}, \phi, \theta | X_U) l(X_m, \bar{U}_m, g_\theta(e_k)) \mathbb{1}_{k=\bar{J}_m} \\
& \quad - \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_m, U_m)} q(\mathbf{e}, \phi, \theta | X_U) l(X_m, U_m, g_\theta(e_k)) \mathbb{1}_{k=J_m} \\
& = \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\bar{J}_m | \mathbf{e}, \phi, X_m, \bar{U}_m)} q(\mathbf{e}, \phi, \theta | X_U) \|X_m, \bar{U}_m - g_\theta(e_k)\|^2 \mathbb{1}_{k=\bar{J}_m} \\
& \quad - \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_m, U_m)} q(\mathbf{e}, \phi, \theta | X_U) \|X_m, U_m - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m}, \quad (4)
\end{aligned}$$

380 where the first term corresponds to the test loss and the second term corresponds to the training loss.

381 Recall the learning algorithm and posterior distribution:

$$\begin{aligned}
\mathbf{e}, \phi, \theta & \sim q(\mathbf{e}, \phi, \theta | X_U), \\
j_k & \sim q(\mathbf{J} | \mathbf{e}, \phi, x_k).
\end{aligned}$$

382 Here $\mathbf{e} = \{e_1, \dots, e_K\}$ is the codebook, and j and $\mathbf{J} = \{J_1, \dots, J_n\}$ represents the index of the
383 codebook that the test and training data are represented.

384 Conditioned on X and U , we then decompose Eq. (4) as follows

$$\begin{aligned}
& \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} l(X_m, \bar{U}_m, g_\theta(e_k)) \mathbb{E}_{q(\bar{J}_m | \mathbf{e}, \phi, X_m, \bar{U}_m)} \mathbb{1}_{k=\bar{J}_m} \\
& \quad - \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} l(X_m, \bar{U}_m, g_\theta(e_k)) \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_m, U_m)} \mathbb{1}_{k=J_m} \\
& \quad + \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} l(X_m, \bar{U}_m, g_\theta(e_k)) \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_m, U_m)} \mathbb{1}_{k=J_m} \\
& \quad - \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} l(X_m, U_m, g_\theta(e_k)) \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_m, U_m)} \mathbb{1}_{k=J_m}. \quad (5)
\end{aligned}$$

385 We will separately upper bound these terms.

386 C.1 Bounding first and second terms

387 The decomposition of the generalization error, as shown in Eq. (5), allows us to bound the first and
388 second terms as follows.

389 We apply Donsker-Varadhan's inequality between the following two distributions:

$$\begin{aligned}
\mathbf{Q} & := P(U) q(\mathbf{e}, \phi, \theta | X_U) q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\bar{U}}, X_U) \\
\mathbf{P}_S & := P(U) q(\mathbf{e}, \phi, \theta | X_U) \mathbb{E}_{P(U')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\bar{U}'}, X_{U'}). \quad (6)
\end{aligned}$$

390 Then, for any $\lambda \in \mathbb{R}^+$, we have

$$\begin{aligned}
& \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} l(X_m, \bar{U}_m, g_\theta(e_k)) \left(\mathbb{E}_{q(\bar{J}_m | \mathbf{e}, \phi, X_m, \bar{U}_m)} \mathbb{1}_{k=\bar{J}_m} - \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_m, U_m)} \mathbb{1}_{k=J_m} \right) \\
& \leq \frac{1}{\lambda} \text{KL}(\mathbf{Q} | \mathbf{P}_S) + \frac{1}{\lambda} \log \mathbb{E}_{\mathbf{P}_S} \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_m, \bar{U}_m, g_\theta(e_k)) (\mathbb{1}_{k=\bar{J}_m} - \mathbb{1}_{k=J_m}) \right).
\end{aligned}$$

391 To simplify the notation, we express $\bar{\mathbf{J}} = \mathbf{J}_0$, $\bar{J}_m = J_{m,0}$, $\mathbf{J} = \mathbf{J}_1$, and $J_m = J_{m,1}$. Let U'' be a
 392 random variable taking 0, 1 with a uniform distribution. Since \mathbf{P}_S is symmetric with respect to the
 393 permutation of \mathbf{J}_0 and \mathbf{J}_1 , we can bound the exponential moment as:

$$\begin{aligned} & \log \mathbb{E}_{P(U)q(\mathbf{e}, \phi, \theta | X_U)} \mathbb{E}_{P(U')} q(\mathbf{J}_0, \mathbf{J}_1 | \mathbf{e}, \phi, X_{\bar{U}'}, X_{U'}) \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{m, \bar{U}_m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{m,0}} - \mathbb{1}_{k=J_{m,1}}) \right) \\ &= \log \mathbb{E}_{P(U)q(\mathbf{e}, \phi, \theta | X_U) P(U'')^n} \mathbb{E}_{P(U')} q(\mathbf{J}_0, \mathbf{J}_1 | \mathbf{e}, \phi, X_{\bar{U}'}, X_{U'}) P(U'')^n \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{m, \bar{U}_m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{m, \bar{U}''}} - \mathbb{1}_{k=J_{m, U''}}) \right) \\ &= \log \mathbb{E}_{P(U)q(\mathbf{e}, \phi, \theta | X_U)} \mathbb{E}_{P(U')} q(\mathbf{J}_0, \mathbf{J}_1 | \mathbf{e}, \phi, X_{\bar{U}'}, X_{U'}) \mathbb{E}_{P(U'')^n} \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{m, \bar{U}_m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{m, \bar{U}''}} - \mathbb{1}_{k=J_{m, U''}}) \right). \end{aligned}$$

394 In the final line, we apply McDiarmid's inequality since U''^n are n i.i.d random variables. To use
 395 McDiarmid's inequality in Lemma 1, we use the stability caused by replacing one of the elements of n
 396 i.i.d random variables. To estimate the coefficients of stability in Lemma 1, let $U''^n = (U''_1, \dots, U''_N)$,
 397 then

$$\begin{aligned} & \sup_{\{U''_m\}_{m=1}^n, U''_{m'}} \left| \frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{m, \bar{U}_m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{m, \bar{U}''_m}} - \mathbb{1}_{k=J_{m, U''_m}}) \right. \\ & \quad \left. - \frac{\lambda}{n} \sum_{k=1}^K \sum_{m \neq m'} l(X_{m, \bar{U}_m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{m, \bar{U}''_m}} - \mathbb{1}_{k=J_{m, U''_m}}) \right. \\ & \quad \left. - \frac{\lambda}{n} \sum_{k=1}^K l(X_{m', \bar{U}'_m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{m', \bar{U}''_{m'}}} - \mathbb{1}_{k=J_{m', U''_{m'}}}) \right| \\ &= \sup_{\{U''_m\}_{m=1}^n, U''_{m'}} \left| \frac{\lambda}{n} \sum_{k=1}^K l(X_{m', \bar{U}'_m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{m', \bar{U}''_{m'}}} - \mathbb{1}_{k=J_{m', U''_{m'}}}) \right. \\ & \quad \left. - \frac{\lambda}{n} \sum_{k=1}^K l(X_{m', \bar{U}'_m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{m', \bar{U}''_{m'}}} - \mathbb{1}_{k=J_{m', U''_{m'}}}) \right| \leq \frac{2\lambda\Delta}{n} \end{aligned} \quad (7)$$

398 Here, the maximum change caused by replacing one element of U'' is $2\lambda\Delta/n$, thus, its log of the
 399 exponential moment is bounded by $(2\lambda\Delta/n)^2/8 \times n = \lambda^2\Delta^2/2n$. Thus from Lemma 1, we have

$$\begin{aligned} & \log \mathbb{E}_{P(U)q(\mathbf{e}, \phi, \theta | X_U)} \mathbb{E}_{P(U')} q(\mathbf{J}_0, \mathbf{J}_1 | \mathbf{e}, \phi, X_{\bar{U}'}, X_{U'}) \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{m, \bar{U}_m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{m,0}} - \mathbb{1}_{k=J_{m,1}}) \right) \\ & \leq \frac{\lambda^2\Delta^2}{2n}. \end{aligned}$$

400 Finally, by noting that

$$\mathbb{E}_X \text{KL}(\mathbf{Q} | \mathbf{P}_S) = \mathbb{E}_X \mathbb{E}_{P(U)q(\mathbf{e}, \phi | X_U)} \text{KL}(q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\bar{U}}, X_U) | \mathbb{E}_{U'} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\bar{U}'}, X_{U'})) = I(\bar{\mathbf{J}}, \mathbf{J}; U | \mathbf{e}, \phi, X),$$

401 the first and second terms in Eq. (5) are upper bounded by

$$\frac{1}{\lambda} I(\bar{\mathbf{J}}, \mathbf{J}; U | \mathbf{e}, \phi, X) + \frac{\lambda\Delta^2}{2n}. \quad (8)$$

402 **C.2 Bounding third and fourth terms**

403 Next, we upper bound the third and fourth terms in Eq.(5);

$$\begin{aligned} & \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} l(X_m, \bar{U}_m, g_\theta(e_k)) \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_m, U_m)} \mathbb{1}_{k=J_m} \\ & - \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} l(X_m, U_m, g_\theta(e_k)) \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_m, U_m)} \mathbb{1}_{k=J_m}. \end{aligned} \quad (9)$$

404 We simplify the notation by expressing $\mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_m, U_m)} \mathbb{1}_{k=J_m}$ as $P_{k,m}$ and use the square loss:

$$\begin{aligned} & \mathbb{E}_{X,U} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} l(X_m, \bar{U}_m, g_\theta(e_k)) P_{k,m} - \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} l(X_m, U_m, g_\theta(e_k)) P_{k,m} \\ & = \mathbb{E}_{X,U} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} (\|X_m, \bar{U}_m\|^2 - \|X_m, U_m\|^2) P_{k,m} \\ & + \mathbb{E}_{X,U} \sum_{k=1}^K \frac{2}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} (X_m, \bar{U}_m - X_m, U_m) \cdot g_\theta(e_k) P_{k,m} \\ & = \mathbb{E}_{X,U} \frac{1}{n} \sum_{m=1}^n (\|X_m, \bar{U}_m\|^2 - \|X_m, U_m\|^2) \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_U)} \sum_{k=1}^K P_{k,m} \\ & + \mathbb{E}_S \frac{2}{n} \sum_{m=1}^n (\mathbb{E}_X X - X_m) \cdot \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \sum_{k=1}^K g_\theta(e_k) P_{k,m} \\ & = \mathbb{E}_S \frac{2}{n} \sum_{m=1}^n (\mathbb{E}_X X - X_m) \cdot \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \sum_{k=1}^K g_\theta(e_k) P_{k,m}, \end{aligned} \quad (10)$$

405 where we express $S = (X_{1,U_1}, \dots, X_{n,U_n}) = (S_1, \dots, S_n)$ as the training samples. In the last
406 inequality, we used $\sum_{k=1}^K P_{k,m} = 1$ and $\mathbb{E}_{X,U} \frac{1}{n} \sum_{m=1}^n (\|X_m, \bar{U}_m\|^2 - \|X_m, U_m\|^2) = 0$ since X
407 and U are i.i.d.

408 To evaluate the final line, we use the Donsker-Valadhan inequality between

$$\begin{aligned} \mathbf{Q} & := q(\mathbf{e}, \phi, \theta | S) \prod_{m=1}^n q(J_m | \mathbf{e}, \phi, S_m), \\ \mathbf{P}_S & := q(\mathbf{e}, \phi, \theta | S) \prod_{m=1}^n q(J_m | \mathbf{e}, \phi), \end{aligned}$$

409 where $q(J_m | \mathbf{e}, \phi)$ is the prior distribution, which never depends on the training data. Then we have

$$\begin{aligned} & \mathbb{E}_S \frac{2}{n} \sum_{m=1}^n (\mathbb{E}_X X - X_m) \cdot \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \sum_{k=1}^K g_\theta(e_k) P_{k,m} \\ & \leq \mathbb{E}_S \frac{1}{\lambda} \text{KL}(\mathbf{Q} | \mathbf{P}_S) + \mathbb{E}_S \frac{1}{\lambda} \log \mathbb{E}_{\mathbf{P}_S} \exp \left(\frac{2\lambda}{n} \sum_{m=1}^n (\mathbb{E}_X X - X_m) \cdot \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \sum_{k=1}^K g_\theta(e_k) \mathbb{1}_{k=J_m} \right) \\ & \leq \mathbb{E}_S \frac{1}{\lambda} \text{KL}(\mathbf{Q} | \mathbf{P}_S) \\ & + \mathbb{E}_S \frac{1}{\lambda} \log \mathbb{E}_{\mathbf{P}_S} \exp \left(\frac{2\lambda}{n} \sum_{m=1}^n (\mathbb{E}_X X - X_m) \cdot \sum_{k=1}^K g_\theta(e_k) (\mathbb{1}_{k=J_m} - P''_{k,m}) \right) \\ & + \mathbb{E}_S \mathbb{E}_{\mathbf{P}_S} \frac{2}{n} \sum_{m=1}^n (\mathbb{E}_X X - X_m) \cdot \sum_{k=1}^K g_\theta(e_k) P''_{k,m}, \end{aligned} \quad (11)$$

410 where $P''_{k,m} = \mathbb{E}_{q(J_m|\phi, \mathbf{e})} \mathbb{1}_{k=J_m}$. Clearly, this does not depend on the index m , so we express
 411 $P''_{k,m} = P''_k$. Then the last term becomes

$$\begin{aligned}
 \mathbb{E}_S \mathbb{E}_{\mathbf{P}_S} \frac{1}{n} \sum_{m=1}^n (\mathbb{E}_X X - X_m) \cdot \sum_{k=1}^K g_\theta(e_k) P''_k &\leq \mathbb{E}_S \mathbb{E}_{\mathbf{P}_S} \left\| \mathbb{E}_X X - \frac{1}{n} \sum_{m=1}^n X_m \right\| \left\| \sum_{k=1}^K g_\theta(e_k) P''_k \right\| \\
 &\leq \mathbb{E}_S \left\| \mathbb{E}_X X - \frac{1}{n} \sum_{m=1}^n X_m \right\| \sqrt{\Delta} \\
 &\leq \sqrt{\Delta \text{Var} \left(\frac{1}{n} \sum_{m=1}^n X_m \right)} \\
 &\leq \sqrt{\Delta \frac{\text{Var}(X)}{n}} \\
 &\leq \sqrt{\frac{\Delta}{4n}} \sqrt{\Delta} = \frac{\Delta}{2\sqrt{n}}, \tag{12}
 \end{aligned}$$

412 where we used the fact that the variance of random variables with bounded in $(a, b]$ is upper bounded
 413 by $(b - a)^2/4n$ (the extension to the d -dimensional random variable is straightforward) and thus,
 414 $\text{Var}(X) \leq \Delta/4$. Then the exponential moment term becomes

$$\begin{aligned}
 &\mathbb{E}_S \frac{1}{\lambda} \log \mathbb{E}_{\mathbf{P}_S} \exp \left(\frac{2\lambda}{n} \sum_{m=1}^n (\mathbb{E}_X X - X_m) \cdot \sum_{k=1}^K g_\theta(e_k) (\mathbb{1}_{k=J_m} - P''_{k,m}) \right) \\
 &= \mathbb{E}_S \frac{1}{\lambda} \log \mathbb{E}_{\mathbf{P}_S} \exp \left(\frac{2\lambda}{n} \sum_{m=1}^n (\mathbb{E}_X X - X_m) \cdot \sum_{k=1}^K g_\theta(e_k) (\mathbb{1}_{k=J} - P''_k) \right).
 \end{aligned}$$

415 Here we use the McDiarmid's inequality for n random variables \mathbf{J} . Then we estimate the stability
 416 coefficient similarly to Eq. (7), which is upper bounded by $\lambda\Delta/n$. Then from Lemma 1, the
 417 exponential moment is bounded by $(2\lambda\Delta/n)^2/8 \times n = \lambda\Delta^2/2n$. Thus, the second term is upper
 418 bounded by

$$\frac{1}{\lambda} \text{KL}(\mathbf{Q}|\mathbf{P}_S) + \frac{\lambda\Delta^2}{2n} + \frac{\Delta}{\sqrt{n}}. \tag{13}$$

419 By optimizing the first and second terms of Eqs. (8) and (13), we have

$$2\Delta \sqrt{\frac{(I(\bar{\mathbf{J}}, \mathbf{J}; U|\mathbf{e}, \phi, X) + \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta|S)} \text{KL}(\mathbf{Q}|\mathbf{P}_S))}{n}} + \frac{\Delta}{\sqrt{n}},$$

420 where we used the fact that X_m are i.i.d. Thus, we use McDiarmid's inequality for n random variables
 421 of X_m to upper bound the exponential moment. We estimate the stability coefficient similarly to
 422 Eq. (7), which is upper bounded by as follows. where

$$\begin{aligned}
 \mathbf{Q} &:= q(\mathbf{e}, \phi, \theta|S) \prod_{m=1}^n q(J_m|\mathbf{e}, \phi, S_m), \\
 \mathbf{P}_S &:= q(\mathbf{e}, \phi, \theta|S) \prod_{m=1}^n q(J_m|\mathbf{e}, \phi).
 \end{aligned}$$

423 C.3 Discussion about the limitation of the existing supersample setting

424 The empirical KL divergence in Theorem 2 originates from the third and fourth terms of Eq.(5), as
 425 discussed in Appendix C.2. After applying the Donsker-Valadhan lemma in the proof, it is crucial to
 426 ensure that the probability $P''_{k,m}$ does not depend on the sample index m to control the exponential
 427 moment in Eq.(11). To achieve this, we employ the prior distribution $q(J_m|\mathbf{e}, \phi)$, which eliminates

428 the sample index dependency and leads to $P''_{k,m} = P''_k$. As a result, we can use a distribution of the
 429 form:

$$\mathbf{P}_S := q(\mathbf{e}, \phi, \theta|S) \prod_{m=1}^n \sum_{m'=1}^n \frac{1}{N} q(J_m|\mathbf{e}, \phi, S_{m'}),$$

430 which provides an empirical approximation of the marginal distribution using available samples.
 431 Since this distribution does not explicitly depend on the sample index, we can bound the exponential
 432 moment similarly as done in Appendix C.2.

433 However, using the prior distribution in Eq.(6) to bound the third and fourth terms of Eq.(5) is not
 434 feasible. The reason is that applying the Donsker-Valadhan lemma with Eq.(6) to these terms does
 435 not yield a bound of order $\mathcal{O}(1/\sqrt{n})$ as achieved in Eq.(12). This is because the dependency on
 436 the sample index in Eq.(6) prevents us from leveraging the symmetry between the test and training
 437 datasets through the supersample index U . Consequently, the prior distribution's symmetry cannot be
 438 exploited to simplify the bounds for these terms.

439 D Proof of Lemmas and equations

440 D.1 Proof of Eq. (2)

441 We define $q(\bar{\mathbf{J}}|\mathbf{e}, \phi) = \prod_{m=1}^n q(\bar{J}_m|\mathbf{e}, \phi)$, $q(\mathbf{J}|\mathbf{e}, \phi) = \prod_{m=1}^n q(J_m|\mathbf{e}, \phi)$, and $q(\tilde{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}) =$
 442 $q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \phi) = q(\bar{\mathbf{J}}|\mathbf{e}, \phi)q(\mathbf{J}|\mathbf{e}, \phi)$ where each $q(\bar{J}_m|\mathbf{e}, \phi)$ is the marginal distribution of
 443 $q(J_m|\mathbf{e}, \phi, X_m)$.

444 Then by the definition of the CMI, we have

$$\begin{aligned} & I(\tilde{\mathbf{J}}; U|\mathbf{e}, \phi, \tilde{X}) \\ &= \mathbb{E}_{\tilde{X}, U} \mathbb{E}_{q(\mathbf{e}, \phi|\tilde{X}_U)} \text{KL}(q(\tilde{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}) \|\mathbb{E}_{U'} q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \phi, \tilde{X}_{\bar{U}'}, \tilde{X}_{U'})) \\ &\leq \mathbb{E}_{\tilde{X}, U} \mathbb{E}_{q(\mathbf{e}, \phi|\tilde{X}_U)} \text{KL}(q(\tilde{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}) \| q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \phi)) \\ &= \mathbb{E}_{\tilde{X}, U} \mathbb{E}_{q(\mathbf{e}, \phi|\tilde{X}_U)} \text{KL}(q(\bar{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}_{\bar{U}}) \| q(\bar{\mathbf{J}}|\mathbf{e}, \phi)) + \mathbb{E}_{\tilde{X}, U} \mathbb{E}_{q(\mathbf{e}, \phi|\tilde{X}_U)} \text{KL}(q(\mathbf{J}|\mathbf{e}, \phi, \tilde{X}_U) \| q(\mathbf{J}|\mathbf{e}, \phi)) \\ &= \mathbb{E}_{\tilde{X}, U} \mathbb{E}_{q(\mathbf{e}, \phi|\tilde{X}_U)} \sum_{m=1}^n \text{KL}(q(\bar{J}_m|\mathbf{e}, \phi, \tilde{X}_{m, \bar{U}_m}) \| q(\bar{J}_m|\mathbf{e}, \phi)) \\ &\quad + \mathbb{E}_{\tilde{X}, U} \mathbb{E}_{q(\mathbf{e}, \phi|\tilde{X}_U)} \sum_{m=1}^n \text{KL}(q(J_m|\mathbf{e}, \phi, \tilde{X}_{m, U_m}) \| q(J_m|\mathbf{e}, \phi)) \\ &= nI(J; X|\mathbf{e}, \phi) + \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi|S)} \frac{1}{n} \sum_{m=1}^n \text{KL}(q(J_m|\mathbf{e}, \phi, S_m) \| q(J_m|\mathbf{e}, \phi)) \\ &\leq nI(e_J; X|\mathbf{e}, \phi) + \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi|S)} \frac{1}{n} \sum_{m=1}^n \text{KL}(q(J_m|\mathbf{e}, \phi, S_m) \| q(J_m|\mathbf{e}, \phi)). \end{aligned}$$

445 D.2 Discussion about the CMI of the deterministic encoder

446 Here, we consider the case where $f_\phi : \mathcal{X} \rightarrow [K]$ represents a deterministic encoder that maps
 447 input data to one of the K indices. This scenario can be interpreted as a K -class classification
 448 problem, allowing us to directly apply the results from Harutyunyan et al. (2021). In their work,
 449 they demonstrated that the CMI for multi-class classification problems can be upper-bounded using
 450 the Natarajan dimension. The Natarajan dimension is a combinatorial measure that generalizes the
 451 VC dimension to multiclass classification setting. Using this concept, we can derive the following
 452 characterization:

453 When using a deterministic encoder network $f_\phi : \mathcal{X} \rightarrow [K]$, belonging to a class with finite Natarajan
 454 dimension d_K , and assuming $2n > d_K + 1$, we have the following bound:

$$I(\tilde{\mathbf{J}}; U|\mathbf{e}, \phi, \tilde{X}) \leq d_K \log \left(\binom{K}{2} \frac{2en}{d_K} \right).$$

455 The proof follows exactly as in Theorem 8 of Harutyunyan et al. (2021).

456 Thus, by regularizing the capacity of the encoder model (via the Natarajan dimension), the CMI term
 457 scales as $\mathcal{O}(\log n)$, ensuring controlled generalization behavior. Examples of models that satisfy the
 458 finite Natarajan dimension are shown in Jin (2023) and Daniely et al. (2011). Also, see Bendavid
 459 et al. (1995), which shows that the VC dimension of the multiclass loss function characterizes the
 460 graph dimension, and the graph dimension upper bounds the Natarajan dimension. For the discussion
 461 of the stochastic encoder that uses $q(J|\mathbf{e}, \phi, x) = q(J|\mathbf{e}, f_\phi(x))$, see Appendix F.2.

462 D.3 Comparison with the fCMI

463 Here, we examine the relationship between our CMI and existing forms of fCMI in more detail. As
 464 highlighted in the main paper, a key difference is that our CMI is conditioned on all model parameters,
 465 whereas existing fCMI approaches marginalize the parameters.

466 To explore this further, we consider marginalizing over the encoder parameter, ϕ . In the proof of
 467 Theorem 2, we perform this marginalization over ϕ in Eq. (4), and obtain

$$\begin{aligned} & \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\bar{J}_m|\mathbf{e}, \phi, X_m, \bar{U}_m)} q(\mathbf{e}, \phi, \theta|X_U) l(X_m, \bar{U}_m, g_\theta(e_k)) \mathbb{1}_{k=\bar{J}_m} \\ & \quad - \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m|\mathbf{e}, \phi, X_m, U_m)} q(\mathbf{e}, \phi, \theta|X_U) l((X_m, U_m, g_\theta(e_k)) \mathbb{1}_{k=J_m} \\ & = \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\bar{J}_m|\theta, \mathbf{e}, X_m, \bar{U}_m)} q(\mathbf{e}, \theta|X_U) \|X_m, \bar{U}_m - g_\theta(e_k)\|^2 \mathbb{1}_{k=\bar{J}_m} \\ & \quad - \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m|\theta, \mathbf{e}, X_m, U_m)} q(\mathbf{e}, \theta|X_U) \|X_m, U_m - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m}, \end{aligned}$$

468 and proceed with the proof in the same way. We apply the Donsker-Varadhan inequality between the
 469 following distributions, instead of Eq.(6):

$$\begin{aligned} \mathbf{Q} & := P(U)P(U')q(\mathbf{e}, \theta|X_U)q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \theta, X_{\bar{U}}, X_U) \\ \mathbf{P} & := P(U)q(\mathbf{e}, \theta|X_U)\mathbb{E}_{P(U')}q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \theta, X_{\bar{U}'}, X_{U'}). \end{aligned}$$

470 This incorporates marginalization over ϕ in Eq.(6), resulting in the following KL divergence in the
 471 upper bound:

$$\begin{aligned} \mathbb{E}_X \text{KL}(\mathbf{Q}|\mathbf{P}) & = \mathbb{E}_{P(U)q(\mathbf{e}, \phi|X_U)} \mathbb{E} \text{KL}(q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \theta, X_{\bar{U}}, X_U)|\mathbb{E}_{P(U')}q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \theta, X_{\bar{U}'}, X_{U'})) \\ & = I(\bar{\mathbf{J}}, \mathbf{J}; U|\mathbf{e}, \theta, X). \end{aligned}$$

472 Unlike Theorem 2, this CMI explicitly involves the decoder parameter θ . By marginalizing over ϕ ,
 473 decoder information is integrated into the upper bound, making Theorem 2 distinct from existing
 474 fCMI bounds.

475 E Proof of Theorem 3

476 We define $\mathbf{T} = \{\mathbf{T}_0, \mathbf{T}_1\}$, where $\tilde{X}_{\mathbf{T}_0} = (\tilde{X}_{T_1}, \dots, \tilde{X}_{T_n})$ serves as the test dataset and $\tilde{X}_{\mathbf{T}_1} =$
 477 $(\tilde{X}_{T_{n+1}}, \dots, \tilde{X}_{T_{2n}})$ serves as the training dataset. We further express $\tilde{X}_{\mathbf{T}_0} = (\tilde{X}_{T_1}, \dots, \tilde{X}_{T_n}) =$
 478 $(\tilde{X}_{\mathbf{T}_0,1}, \dots, \tilde{X}_{\mathbf{T}_0,n})$ and $\tilde{X}_{\mathbf{T}_1} = (\tilde{X}_{\mathbf{T}_1,1}, \dots, \tilde{X}_{\mathbf{T}_1,n})$. To emphasize the dependence of the
 479 dataset on \mathbf{T} , we write the posterior distribution as $q(\bar{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}}) = q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}}) =$
 480 $q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}_0}, \tilde{X}_{\mathbf{T}_1}) = q(\bar{\mathbf{J}}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}_0})q(\mathbf{J}|\mathbf{e}, \phi, \tilde{X}_{\mathbf{T}_1})$.

481 Hereinafter, we express \tilde{X} as X to simplify the notation. Under the permutation symmetric settings,
 482 the generalization error can be expressed as

$$\begin{aligned}
& \mathbb{E}_{S, X} \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \left(\mathbb{E}_{q(J | \mathbf{e}, \phi, X)} l(X, g_\theta(e_J)) - \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \phi, S_m)} l(S_m, g_\theta(e_{J_m})) \right) \\
&= \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\bar{J}_m | \mathbf{e}, \phi, X_{\mathbf{T}_0, m}) q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} l((X_{\mathbf{T}_0, m}, g_\theta(e_k)) \mathbb{1}_{k=\bar{J}_m} \\
&- \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_{\mathbf{T}_1, m}) q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} l(X_{\mathbf{T}_1, m}, g_\theta(e_k)) \mathbb{1}_{k=J_m} \\
&= \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\bar{J}_m | \mathbf{e}, \phi, X_{\mathbf{T}_0, m}) q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_0, m} - g_\theta(e_k)\|^2 \mathbb{1}_{k=\bar{J}_m} \\
&- \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_{\mathbf{T}_1, m}) q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_1, m} - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m}. \tag{14}
\end{aligned}$$

483 We then decompose the loss as follows

$$\begin{aligned}
& \text{gen}(n, \mathcal{D}) \tag{15} \\
&= \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\bar{J}_m | \mathbf{e}, \phi, X_{\mathbf{T}_0, m}) q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_0, m} - g_\theta(e_k)\|^2 \mathbb{1}_{k=\bar{J}_m} \\
&- \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_{\mathbf{T}_1, m}) q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_0, m} - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m} \\
&+ \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_{\mathbf{T}_1, m}) q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_0, m} - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m} \\
&- \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_{\mathbf{T}_1, m}) q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_1, m} - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m}.
\end{aligned}$$

484 First, we upper bound the first two terms by applying the Donsker-Varadhan inequality. Consider the
 485 joint distribution and the prior distribution, defined as follows:

$$\begin{aligned}
\mathbf{Q} &:= P(\mathbf{T}) q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1}) q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}}), \\
\mathbf{P} &:= P(\mathbf{T}) q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1}) \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}}).
\end{aligned}$$

486 Then we then obtain

$$\begin{aligned}
& \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_0, m} - g_\theta(e_k)\|^2 \left(\mathbb{E}_{q(\bar{J}_m | \mathbf{e}, \phi, X_{\mathbf{T}_1, m})} \mathbb{1}_{k=\bar{J}_m} - \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_{\mathbf{T}_0, m})} \mathbb{1}_{k=J_m} \right) \\
&\leq \mathbb{E}_X \frac{1}{\lambda} \text{KL}(\mathbf{Q} | \mathbf{P}) + \mathbb{E}_X \frac{1}{\lambda} \log \mathbb{E}_{\mathbf{P}} \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n \|X_{\mathbf{T}_0, m} - g_\theta(e_k)\|^2 (\mathbb{1}_{k=\bar{J}_m} - \mathbb{1}_{k=J_m}) \right).
\end{aligned}$$

487 Note that $\mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'})$ is symmetric with respect to the permutation of \mathbf{T} . Thus, we have

$$\begin{aligned}
& \log \mathbb{E}_{P(\mathbf{T}) q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'}) \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{\mathbf{T}_0, m}, g_\theta(e_k)) (\mathbb{1}_{k=\bar{J}_m} - \mathbb{1}_{k=J_m}) \right) \\
&= \log \mathbb{E}_{P(\mathbf{T}) q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'}) \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{\mathbf{T}_0, m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{\mathbf{T}''_0, m}} - \mathbb{1}_{k=J_{\mathbf{T}''_1, m}}) \right) \\
&= \log \mathbb{E}_{P(\mathbf{T}) q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'}) \mathbb{E}_{P(\mathbf{T}'')} \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{\mathbf{T}_0, m}, g_\theta(e_k)) (\mathbb{1}_{k=J_{\mathbf{T}''_0, m}} - \mathbb{1}_{k=J_{\mathbf{T}''_1, m}}) \right).
\end{aligned}$$

488 To simplify the notation, we define $\mathbf{T}'' = \{\mathbf{T}''_0, \mathbf{T}''_1\} = \{\mathbf{T}''_{0,1}, \dots, \mathbf{T}''_{0,n}, \mathbf{T}''_{1,1}, \dots, \mathbf{T}''_{1,n}\}$. Note
489 that $\mathbf{T}''_{j,m}$ for $m = 1, \dots, n$ and $j = 0, 1$ are not independent of each other due to the permutation
490 that generates them. Therefore, we cannot directly apply standard concentration inequalities, as is
491 possible in the existing supersample setting.

492 To address this, we use the results from Joag-Dev & Proschan (1983), which concern the negative
493 association of permutation variables. From Theorem 2.11 in Joag-Dev & Proschan (1983), the
494 distribution $P(\mathbf{T})$ satisfies negative association. Additionally, as discussed in Section 3.3 of Joag-Dev
495 & Proschan (1983) and further in Proposition 4 and 5 of Dubhashi & Ranjan (1996), we have that

$$\begin{aligned} & \log \mathbb{E}_{P(\mathbf{T})q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'}) \mathbb{E}_{P(\mathbf{T}'')} \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{\mathbf{T}_{0,m}}, g_\theta(e_k)) (\mathbb{1}_{k=J_{\mathbf{T}''_{0,m}}} - \mathbb{1}_{k=J_{\mathbf{T}''_{1,m}}}) \right) \\ & \leq \log \mathbb{E}_{P(\mathbf{T})q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'}) \mathbb{E}_{\prod_{m=1}^n \prod_{j=0,1} P(\mathbf{T}''_{j,m})} \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{\mathbf{T}_{0,m}}, g_\theta(e_k)) (\mathbb{1}_{k=J_{\mathbf{T}''_{0,m}}} - \mathbb{1}_{k=J_{\mathbf{T}''_{1,m}}}) \right) \end{aligned}$$

496 where $P(\mathbf{T}''_{j,m})$ is the marginal distribution, implying that $\mathbf{T}''_{j,m}$ are now $2n$ independent random
497 variables. Intuitively, the results in Joag-Dev & Proschan (1983) indicate that the elements of the
498 permutation index, which follow the permutation distribution, are negatively correlated. As a result,
499 the expectation of the marginal distribution is larger than that of the joint distribution.

500 Since $\{\mathbf{T}''_{j,m}\}$ are independent, we can apply McDiarmid's inequality, which leads to the results in

$$\begin{aligned} & \log \mathbb{E}_{P(\mathbf{T})q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'}) \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{\mathbf{T}_{0,m}}, g_\theta(e_k)) (\mathbb{1}_{k=J_m} - \mathbb{1}_{k=J_m}) \right) \\ & \leq \log \mathbb{E}_{P(\mathbf{T})q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'}) \mathbb{E}_{\prod_{m=1}^n \prod_{j=0,1} P(\mathbf{T}''_{j,m})} \exp \left(\frac{\lambda}{n} \sum_{k=1}^K \sum_{m=1}^n l(X_{\mathbf{T}_{0,m}}, g_\theta(e_k)) (\mathbb{1}_{k=J_{\mathbf{T}''_{0,m}}} - \mathbb{1}_{k=J_{\mathbf{T}''_{1,m}}}) \right) \\ & \leq \frac{\lambda^2 \Delta^2}{n}. \end{aligned} \tag{16}$$

501 which is estimated similarly to Eq. (7). Note that there are $2n$ variables so the calculation of the upper
502 bound is $(2\Delta\lambda/n)^2/8 \times 2n = \lambda^2 \Delta^2/n$.

503 Next we focus on the third and fourth terms in Eq. (15). Similarly to Eq. (10), we have

$$\begin{aligned} & \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_{\mathbf{T}_{1,m}})q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_{0,m}} - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m} \\ & - \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_{\mathbf{T}_{1,m}})q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_{1,m}} - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m} \\ & = \mathbb{E}_{X, \mathbf{T}} \frac{2}{n} \sum_{m=1}^n (X_{\mathbf{T}_{1,m}} - X_{\mathbf{T}_{0,m}}) \cdot \mathbb{E}_{q(J_m | \mathbf{e}, \phi, X_{\mathbf{T}_{1,m}})q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \sum_{k=1}^K g_\theta(e_k) \mathbb{1}_{k=J_m} \\ & \leq \mathbb{E}_X \frac{1}{\lambda} \text{KL}(\mathbf{Q} | \mathbf{P}) + \mathbb{E}_X \frac{1}{\lambda} \log \mathbb{E}_{\mathbf{P}} \exp \left(\frac{2\lambda}{n} \sum_{m=1}^n (X_{\mathbf{T}_{1,m}} - X_{\mathbf{T}_{0,m}}) \cdot \sum_{k=1}^K g_\theta(e_k) \mathbb{1}_{k=J_m} \right) \\ & \leq \mathbb{E}_X \frac{1}{\lambda} \text{KL}(\mathbf{Q} | \mathbf{P}) \\ & + \mathbb{E}_X \frac{1}{\lambda} \log \mathbb{E}_{P(\mathbf{T})q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'}) \mathbb{E}_{\prod_{m=1}^n \prod_{j=0,1} P(\mathbf{T}''_{j,m})} \\ & \exp \left(\frac{2\lambda}{n} \sum_{m=1}^n (X_{\mathbf{T}_{1,m}} - X_{\mathbf{T}_{0,m}}) \cdot \sum_{k=1}^K g_\theta(e_k) \mathbb{1}_{k=J_m} \right). \end{aligned} \tag{17}$$

504 We first evaluate the expectation of the exponential moment;

$$\Omega := \mathbb{E}_{P(\mathbf{T})} \mathbb{E}_{q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \frac{2}{n} \sum_{m=1}^n (X_{\mathbf{T}_{1,m}} - X_{\mathbf{T}_{0,m}}) \cdot \mathbb{E}_{P(\mathbf{T}')} \mathbb{E}_{q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'})} \sum_{k=1}^K g_\theta(e_k) \mathbb{1}_{k=J_m}. \quad (18)$$

505 Let us now focus on the expectation $\mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'})$. Due to the permutation symmetry,

506 $\mathbb{E}_{P(\mathbf{T}')} \mathbb{E}_{q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'})} \sum_{k=1}^K g_\theta(e_k) \mathbb{1}_{k=J_m}$ is the same for all m .

507 For instance, when $n = 2$, the possible permutations of \mathbf{T} are $\mathbf{T} =$
508 $(1, 2, 3, 4), (1, 2, 4, 3), (1, 3, 2, 4), \dots$, resulting in 24 distinct patterns and thus

$$\begin{aligned} P_{k,1} &= \mathbb{E}_{P(\mathbf{T}')} \mathbb{E}_{q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'})} \mathbb{1}_{k=\bar{J}_1} = \mathbb{E}_{\frac{1}{4}q(J_1 | \mathbf{e}, \phi, X_1) + \frac{1}{4}q(J_1 | \mathbf{e}, \phi, X_2) + \frac{1}{4}q(J_1 | \mathbf{e}, \phi, X_3) + \frac{1}{4}q(J_1 | \mathbf{e}, \phi, X_4)} \mathbb{1}_{k=J_1} \\ P_{k,2} &= \mathbb{E}_{P(\mathbf{T}')} \mathbb{E}_{q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'})} \mathbb{1}_{k=\bar{J}_2} = \mathbb{E}_{\frac{1}{4}q(J_2 | \mathbf{e}, \phi, X_1) + \frac{1}{4}q(J_2 | \mathbf{e}, \phi, X_2) + \frac{1}{4}q(J_2 | \mathbf{e}, \phi, X_3) + \frac{1}{4}q(J_2 | \mathbf{e}, \phi, X_4)} \mathbb{1}_{k=J_2} \\ &\vdots \end{aligned}$$

509 Thus, all $P_{k,m}$ does not depend on the index m . So we express

510 $\mathbb{E}_{P(\mathbf{T}')} \mathbb{E}_{q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \phi, X_{\mathbf{T}'})} \sum_{k=1}^K g_\theta(e_k) \mathbb{1}_{k=J_m}$ as P_k . Then Eq. (18) can be written as

$$\begin{aligned} &\mathbb{E}_X \mathbb{E}_{P(\mathbf{T})} \mathbb{E}_{q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \left(\frac{1}{n} \sum_{m=1}^n X_{\mathbf{T}_{1,m}} - \frac{1}{n} \sum_{m=1}^n X_{\mathbf{T}_{0,m}} \right) \cdot \sum_{k=1}^K g_\theta(e_k) P_k \\ &\leq \mathbb{E}_X \mathbb{E}_{P(\mathbf{T})} \mathbb{E}_{q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \left\| \frac{1}{n} \sum_{m=1}^n X_{\mathbf{T}_{1,m}} - \frac{1}{n} \sum_{m=1}^n X_{\mathbf{T}_{0,m}} \right\| \left\| \mathbb{E}_{P(\mathbf{T})} \mathbb{E}_{q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \left\| \sum_{k=1}^K g_\theta(e_k) P_k \right\| \right\| \\ &\leq \mathbb{E}_X \mathbb{E}_{P(\mathbf{T})} \left\| \frac{1}{n} \sum_{m=1}^n X_{\mathbf{T}_{1,m}} - \frac{1}{n} \sum_{m=1}^n X_{\mathbf{T}_{0,m}} \right\| \sqrt{\Delta} \\ &\leq \mathbb{E}_X \mathbb{E}_{\prod_{m=1}^n \prod_{j=0,1} P(\mathbf{T}''_{j,m})} \left\| \frac{1}{n} \sum_{m=1}^n X_{\mathbf{T}_{1,m}} - \frac{1}{n} \sum_{m=1}^n X_{\mathbf{T}_{0,m}} \right\| \sqrt{\Delta}, \end{aligned}$$

511 where we used the negative association property of the permutation distribution. We bound the above
512 exactly same ways as Eq. (12), that is, we can upper bound the above by the variance of bounded
513 random variable and thus, we have

$$\mathbb{E}_X \mathbb{E}_{\prod_{m=1}^n \prod_{j=0,1} P(\mathbf{T}''_{j,m})} \left\| \frac{1}{n} \sum_{m=1}^n X_{\mathbf{T}_{1,m}} - \frac{1}{n} \sum_{m=1}^n X_{\mathbf{T}_{0,m}} \right\| \leq 2\sqrt{\frac{\Delta}{4n}}.$$

514 Thus, we have

$$\Omega = \mathbb{E}_X \mathbb{E}_{P(\mathbf{T})} \mathbb{E}_{q(\mathbf{e}, \theta, \phi | X_{\mathbf{T}_1})} \left(\frac{2}{n} \sum_{m=1}^n X_{\mathbf{T}_{1,m}} - \frac{2}{n} \sum_{m=1}^n X_{\mathbf{T}_{0,m}} \right) \cdot \sum_{k=1}^K g_\theta(e_k) P_k \leq \frac{2\Delta}{\sqrt{n}},$$

515 Let us back to the evaluation of the exponential moment in Eq. (17), we will evaluate the following

$$\mathbb{E}_X \frac{1}{\lambda} \text{KL}(\mathbf{Q} | \mathbf{P}) + \mathbb{E}_X \frac{1}{\lambda} \log \mathbb{E}_{\mathbf{P}} \exp \left(\frac{2\lambda}{n} \sum_{m=1}^n (X_{\mathbf{T}_{1,m}} - X_{\mathbf{T}_{0,m}}) \cdot \sum_{k=1}^K g_\theta(e_k) \mathbb{1}_{k=J_m} - \lambda \Omega \right) + \Omega.$$

516 We then evaluate this similarly to Eq. (16), which uses the negative association of the permuta-
517 tion distribution and McDiarmid's inequality. The the exponential moment is upper bounded by

518 $(2\Delta\lambda/n)^2/8 \times 2n = \lambda^2\Delta^2/n$ We then obtain

$$\begin{aligned}
& \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m|\mathbf{e}, \phi, X_{\mathbf{T}_{1,m}})q(\mathbf{e}, \phi, \theta|X_{\mathbf{T}_1})} \|X_{\mathbf{T}_{1,m}} - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m} \\
& - \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m|\mathbf{e}, \phi, X_{\mathbf{T}_{0,m}})q(\mathbf{e}, \phi, \theta|X_{\mathbf{T}_1})} \|X_{\mathbf{T}_{0,m}} - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m} \\
& \leq \mathbb{E}_X \frac{1}{\lambda} \text{KL}(\mathbf{Q}|\mathbf{P}) + \mathbb{E}_X \frac{1}{\lambda} \log \mathbb{E}_{\mathbf{P}} \exp \left(\frac{2\lambda}{n} \sum_{m=1}^n (X_{\mathbf{T}_{1,m}} - X_{\mathbf{T}_{0,m}}) \cdot \sum_{k=1}^K g_\theta(e_k) \mathbb{1}_{k=J_m} - \lambda\Omega \right) + \Omega \\
& \leq \mathbb{E}_X \frac{1}{\lambda} \text{KL}(\mathbf{Q}|\mathbf{P}) + \frac{\lambda\Delta^2}{n} + \frac{2\Delta}{\sqrt{n}}. \tag{19}
\end{aligned}$$

519 In conclusion, from Eqs. (16) and (19) we have

$$\text{gen}(n, \mathcal{D}) \leq \mathbb{E}_X \frac{2}{\lambda} \text{KL}(\mathbf{Q}|\mathbf{P}) + \frac{2\lambda\Delta^2}{n} + \frac{2\Delta}{\sqrt{n}},$$

520 and optimizing the λ , we have

$$\text{gen}(n, \mathcal{D}) \leq 4\Delta \mathbb{E}_X \sqrt{\frac{\text{KL}(\mathbf{Q}|\mathbf{P})}{n}} + \frac{2\Delta}{\sqrt{n}} = 4\Delta \sqrt{\frac{I(\bar{\mathbf{J}}, \mathbf{J}; \mathbf{T}|\mathbf{e}, \phi, X)}{n}} + \frac{2\Delta}{\sqrt{n}}.$$

521 F Proof of Theorem 4

522 Here, we present the results for a general stochastic encoder. For fixed ϕ and \mathbf{e} , assume that for
523 all $\mathbf{x} \in \tilde{X}$, for any $j \in [K]$, and for a fixed $\delta \in \mathbb{R}^+$, the following holds: $q(J = j|\mathbf{e}, f_\phi(x)) \leq$
524 $e^{h(\delta)} q(J = j|\mathbf{e}, \hat{f}(x))$ with $h : \mathbb{R}^+ \rightarrow \mathbb{R}^+$.

525 **Theorem 6.** Assume that there exists a positive constant Δ_z such that $\sup_{z, z' \in \mathcal{Z}} \|z - z'\| < \Delta_z$.
526 Then, when using Eq. (1) and under the same setting as Theorem 3, for any $\delta \in (0, 1]$, we have

$$\text{gen}(n, \mathcal{D}) \leq \Delta \sqrt{nh(\delta)} + 4\Delta \sqrt{\frac{\log \mathcal{N}(\delta, \mathcal{F}, 2n)}{n}} + \frac{2\Delta}{\sqrt{n}}.$$

527 To prove this lemma, we first replace the output of the encoder with that obtained using the δ -cover
528 of the encoder network. Since we assumed that $q(J = j|\mathbf{e}, \phi, x) = q(J = j|\mathbf{e}, f_\phi(x))$, we need to
529 approximate the error caused by $q(J = j|\mathbf{e}, \hat{f}(x))$ approximating $q(J = j|\mathbf{e}, \phi, x)$. To evaluate this
530 gap, we apply the Donsker-Valadhan lemma between the two distributions

$$\begin{aligned}
\mathbf{Q} & := q(J|\mathbf{e}, f_\phi(X)) \prod_{m=1}^n q(J_m|\mathbf{e}, f_\phi(S_m)), \\
\mathbf{P} & := q(J|\mathbf{e}, \hat{f}(X)) \prod_{m=1}^n q(J_m|\mathbf{e}, \hat{f}(S_m)). \tag{20}
\end{aligned}$$

531 Then we have

$$\begin{aligned}
& \text{gen}(n, \mathcal{D}) \\
&= \mathbb{E}_{S, X} \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \left(\mathbb{E}_{q(J | \mathbf{e}, f_\phi(X))} l(X, g_\theta(e_J)) - \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, f_\phi(S_m))} l(S_m, g_\theta(e_{J_m})) \right) \\
&\leq \mathbb{E}_{S, X} \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \left(\mathbb{E}_{q(J | \mathbf{e}, \hat{f}(X))} l(X, g_\theta(e_J)) - \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \hat{f}(S_m))} l(S_m, g_\theta(e_{J_m})) \right) \\
&+ \mathbb{E}_{S, X} \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \frac{1}{\lambda} \text{KL}(\mathbf{Q} \| \mathbf{P}) \\
&+ \mathbb{E}_{S, X} \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \frac{1}{\lambda} \log \mathbb{E}_{\mathbf{P}} \exp \left(\lambda l(X, g_\theta(e_J)) - \frac{\lambda}{n} \sum_{m=1}^n l(S_m, g_\theta(e_{J_m})) \right) \\
&\leq \mathbb{E}_{S, X} \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \left(\mathbb{E}_{q(J | \mathbf{e}, \hat{f}(X))} l(X, g_\theta(e_J)) - \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \hat{f}(S_m))} l(S_m, g_\theta(e_{J_m})) \right) \\
&+ \frac{(n+1)h(\delta)}{\lambda} + \frac{\lambda \Delta^2}{2},
\end{aligned}$$

532 where we used the following relation

$$\text{KL}(\mathbf{Q} \| \mathbf{P}) \leq (n+1) \log e^{h(\delta)} = (n+1)h(\delta),$$

533 which is proved by the assumption of the stability. We also used the fact that $-\lambda \Delta \leq \lambda l(X, g_\theta(e_J)) -$

534 $\frac{\lambda}{n} \sum_{m=1}^n l(S_m, g_\theta(e_{J_m})) \leq \lambda \Delta$ to upper bound the exponential moment.

535 By optimizing λ , we have

$$\begin{aligned}
& \text{gen}(n, \mathcal{D}) \\
&\leq \mathbb{E}_{S, X} \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \left(\mathbb{E}_{q(J | \mathbf{e}, \hat{f}(X))} l(X, g_\theta(e_J)) - \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \hat{f}(S_m))} l(S_m, g_\theta(e_{J_m})) \right) \\
&+ \Delta \sqrt{\frac{(n+1)h(\delta)}{2}}.
\end{aligned}$$

536 This implies that the first term corresponds to the generalization bound when using the δ -cover of the
537 encoder network. We can bound this term similarly to Theorem 3.

538 When applying the result of Theorem 3, we utilize the Donsker-Valadhan inequality for Eq.(14).

539 Instead of using Eq.(20), we consider the following distributions:

$$\begin{aligned}
\mathbf{Q} &:= q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, f_\phi(X_{\mathbf{T}})) = \prod_{m=1}^n q(\bar{J}_m | \mathbf{e}, \hat{f}(\tilde{X}_{T_m})) \prod_{m=1}^n q(J_m | \mathbf{e}, \hat{f}(\tilde{X}_{T_{n+m}})) \\
\mathbf{P} &:= q(\bar{\mathbf{J}}, \mathbf{J} | \mathbf{e}, \hat{f}(X_{\mathbf{T}})) = \prod_{m=1}^n q(\bar{J}_m | \mathbf{e}, \hat{f}(\tilde{X}_{T_m})) \prod_{m=1}^n q(J_m | \mathbf{e}, \hat{f}(\tilde{X}_{T_{n+m}})).
\end{aligned}$$

540 From assumption, we have

$$\text{KL}(\mathbf{Q} \| \mathbf{P}) \leq 2n \log e^{h(\delta)} = 2nh(\delta).$$

541 Then from Eq. (14), we have

$$\begin{aligned}
& \text{gen}(n, \mathcal{D}) \\
&\leq \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(\bar{J}_m | \mathbf{e}, \hat{f}(X_{\mathbf{T}_0, m}))} \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_0, m} - g_\theta(e_k)\|^2 \mathbb{1}_{k=\bar{J}_m} \\
&- \mathbb{E}_{X, \mathbf{T}} \sum_{k=1}^K \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(J_m | \mathbf{e}, \hat{f}(X_{\mathbf{T}_1, m}))} \mathbb{E}_{q(\mathbf{e}, \phi, \theta | X_{\mathbf{T}_1})} \|X_{\mathbf{T}_1, m} - g_\theta(e_k)\|^2 \mathbb{1}_{k=J_m} + \Delta \sqrt{\frac{2nh(\delta)}{2}} \\
&\leq 4\Delta \sqrt{\frac{I(\bar{\mathbf{J}}, \mathbf{J}; \mathbf{T} | \mathbf{e}, \hat{f}(X))}{n}} + \frac{2\Delta}{\sqrt{n}} + \Delta \sqrt{\frac{2nh(\delta)}{2}},
\end{aligned}$$

542 where

$$I(\bar{\mathbf{J}}, \mathbf{J}; \mathbf{T}|\mathbf{e}, \phi, X) = \mathbb{E}_{\tilde{X}, \mathbf{T}} \mathbb{E}_{q(\mathbf{e}, \phi|\tilde{X}_{\mathbf{T}_1})} \text{KL}(q(\bar{\mathbf{J}}|\mathbf{e}, \hat{f}(\tilde{X})) \| \mathbb{E}_{P(\mathbf{T}')} q(\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \hat{f}(\tilde{X}_{\mathbf{T}'_0}), \hat{f}(\tilde{X}_{\mathbf{T}'_1}))).$$

543 Note that we consider the CMI for the discrete variable, it is upper bounded by the entropy (Cover &
544 Thomas, 2012), and we have

$$I(\bar{\mathbf{J}}, \mathbf{J}; \mathbf{T}|\mathbf{e}, \hat{f}(X)) \leq H[\bar{\mathbf{J}}, \mathbf{J}|\mathbf{e}, \hat{f}(X)] \leq \log \mathcal{N}(\delta, \mathcal{F}, 2n).$$

545 The first inequality follows from the fact that MI is defined as the difference between the entropy
546 and the conditional entropy, and the entropy of discrete variables is always non-negative. The
547 second inequality arises because $\bar{\mathbf{J}}, \mathbf{J}$ are outputs of a function evaluated at $2n$ points. Thus, we
548 considered the covering number at $2n$ points, defined as $\mathcal{N}(\delta, \mathcal{F}, n) := \sup_{x^{2n} \in \mathcal{X}^{2n}} \mathcal{N}(\delta, \mathcal{F}, x^{2n})$.
549 Since the entropy is bounded above by the logarithm of the maximum cardinality, we obtain the
550 second inequality.

551 Thus, we have

$$\text{gen}(n, \mathcal{D}) \leq 4\Delta \sqrt{\frac{\log \mathcal{N}(\delta, \mathcal{F}, 2n)}{n}} + \frac{2\Delta}{\sqrt{n}} + \Delta \sqrt{nh(\delta)}.$$

552 F.1 Behavior of Eq. (1)

553 Finally, we show that Eq. (1) satisfies $h(\delta) = 8\beta\Delta_z\delta$ because

$$\begin{aligned} \frac{q(J = j|\mathbf{e}, f_\phi(x))}{q(J = j|\mathbf{e}, \hat{f}(x))} &= \frac{e^{-\beta\|f_\phi(x) - e_j\|^2}}{e^{-\beta\|\hat{f}(x) - e_j\|^2}} \times \frac{\sum_{k=1}^K e^{-\beta\|\hat{f}(x) - e_k\|^2}}{\sum_{k=1}^K e^{-\beta\|f_\phi(x) - e_k\|^2}} \\ &= e^{-\beta\|f_\phi(x) - e_j\|^2 + \beta\|\hat{f}(x) - e_j\|^2} \times \frac{\sum_{k=1}^K e^{\beta\|f_\phi(x) - e_k\|^2}}{\sum_{k=1}^K e^{\beta\|\hat{f}(x) - e_k\|^2}} \\ &\leq e^{\beta(\hat{f}(x) - f_\phi(x)) \cdot (\hat{f}(x) + f_\phi(x)) - 2\beta e_j \cdot (\hat{f}(x) - f_\phi(x))} \times \sup_{k \in [K]} e^{-\beta\|\hat{f}(x) - e_k\|^2 + \beta\|f_\phi(x) - e_k\|^2} \\ &\leq e^{4\beta\Delta_z\delta} \times e^{4\beta\Delta_z\delta}. \end{aligned}$$

554 Thus we have $h(\delta) = 8\beta\Delta_z\delta$ and by substituting this into above Theorem, we obtain Theorem 4.

555 F.2 Discussion about the metric entropy for regularized model

556 Here we discuss the upper bound of metric entropy in our setting. Since the latent variable lies in
557 \mathbb{R}^{d_z} , the encoder network operates as $f_\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{d_z}$, making it a multivariate function.

558 To evaluate the complexity of the metric entropy for such multivariate functions, the concept of
559 Natarajan dimension with margin has been employed (Guermeur, 2007). According to Lemma 39
560 (and also Lemma 37 and 38), if a multivariate function has a finite Natarajan dimension with margin,
561 then its metric entropy scales as $\mathcal{O}(\log n)$. To explore the properties of the Natarajan dimension with
562 margin, Guermeur (2018) demonstrated that it can be bounded by the fat-shattering dimension of
563 each component of the original multivariate function (Lemma 10). Additionally, Guermeur (2017)
564 showed in Lemma 1 that the covering number of the multivariate function can be bounded by the
565 covering number of each of its components. To further bound the covering number of each dimension,
566 one can rely on the fat-shattering dimension of each function, as discussed in Lemma 3.5 of Alon
567 et al. (1997).

568 Thus, it is essential to bound the fat-shattering dimension in both cases. Examples of fat-shattering
569 dimension evaluations can be found, for instance, in Bartlett & Maass (2003), which discusses neural
570 network models, and Gottlieb et al. (2014), which addresses the fat-shattering dimension of Lipschitz
571 function classes. If our encoder network adheres to these properties, we can bound its covering
572 number accordingly.

573 In conclusion, if the log of the covering number satisfies $\mathcal{O}(\log n)$, by setting $\delta = 1/n^2$, we obtain
574 that $\text{gen}(n, \mathcal{D}) = \mathcal{O}(\sqrt{\log n/n})$.

575 G Data generation guarantee for the encoder-decoder model

576 The primary interest of latent variable models often lies in the data generation ability rather than
 577 their generalization under the reconstruction loss. Specifically, the aim is to generate realistic data by
 578 sampling from the latent variable distribution and transforming it via the decoder. We expect that the
 579 distribution of generated data is close to the true data distribution.

580 Let p represent a distribution on the latent space \mathcal{Z} , and assume that for any $\theta \in \Theta$, the decoder
 581 $g_\theta(\cdot) : \mathcal{Z} \rightarrow \mathcal{X}$ is measurable. The pushforward of the distribution p by the decoder, denoted as
 582 $g_\theta\#p$, defines a distribution on \mathcal{X} as $g_\theta\#p(A) = p(g_\theta^{-1}(A))$ for any measurable set $A \subseteq \mathcal{X}$. When
 583 generating data, we first draw a index using prior distribution $p(J|\mathbf{e}, \phi)$, which is typically independent
 584 of the training dataset. This corresponds to selecting a latent variable e_J from $\{e_1, \dots, e_K\}$, and we
 585 denote the associated prior distribution over \mathcal{Z} as $p(e|\mathbf{e}, \phi)$. The resulting distribution of the generated
 586 data is then represented as $\hat{\mu} := g_\theta\#p(e|\mathbf{e}, \phi)$. Next, given the posterior distribution $q(J_m|\mathbf{e}, \phi, S_m)$
 587 conditioned on the m -th training data point S_m , we express the corresponding posterior distribution
 588 over \mathcal{Z} as $q(e_{(m)}|\mathbf{e}, \phi, S_m)$, where we simply express e_{J_m} as $e_{(m)}$. Here, our goal is to bound the
 589 2-Wasserstein distance (See Appendix A for the definition) between data distribution \mathcal{D} and the
 590 data-generating distribution $\hat{\mu}$, denoted as $W_2(\mathcal{D}, \hat{\mu})$. Following is our main result:

591 **Theorem 7.** *Let $S = (S_1, \dots, S_n) \in \mathcal{X}^n$ be a training dataset, where $S_m \in \mathcal{X}$ are drawn i.i.d from*
 592 *\mathcal{D} . Under Assumption 1 and for any prior $q(e|\mathbf{e}, \phi)$ that does not depend on S , we have:*

$$\begin{aligned} \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta|S)} W_2^2(\mathcal{D}, \hat{\mu}) &\leq \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta|S)} \frac{2}{n} \sum_{m=1}^n \mathbb{E}_{q(e_{(m)}|\mathbf{e}, \phi, S_m)} \|S_m - g_\theta(e_{(m)})\|^2 \\ &\quad + 2\Delta \sqrt{\frac{2}{n} \sum_{m=1}^n \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi|S)} \text{KL}(q(e_{(m)}|\mathbf{e}, \phi, S_m) \| q(e_{(m)}|\mathbf{e}, \phi))} + \frac{2\Delta}{\sqrt{n}}. \end{aligned}$$

593 This theorem shows that the 2-Wasserstein distance is upper-bounded by the reconstruction loss on
 594 the training dataset and an empirical KL term. The result is similar to the bound in Mbacke et al.
 595 (2023), which assumes the fixed parameters, that is, learning is not considered. In contrast, our bound
 596 incorporates the learning process of parameters. If the marginal distribution of $q(e|\mathbf{e}, \phi, x)$ were used
 597 as the prior distribution, the empirical KL term would become the empirical MI as discussed in
 598 Sec. 3.1. Furthermore, if a prior distribution with the symmetry introduced in Sec. 3.2 were used,
 599 the empirical KL term would become the CMI appearing in Theorem 3. However, such priors are
 600 impractical in real-world scenarios, where uniform distributions are typically used to sample latent
 601 variables.

602 H Proof of Theorem 7

603 Define the distribution obtained by the training dataset as follows; conditioned on \mathbf{e}, ϕ, S , we have

$$\hat{\mu}_S = \frac{1}{n} \sum_{m=1}^n g_\theta\#q(e_{(m)}|\mathbf{e}, \phi, S_m)$$

604 From the triangle inequality, we have

$$W_2(\mathcal{D}, \hat{\mu}) \leq W_2(\mathcal{D}, \hat{\mu}_S) + W_2(\hat{\mu}_S, \hat{\mu}) \quad (21)$$

605 The first term of Eq. (21) is bounded as follows;

$$\begin{aligned} \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta|S)} W_2^2(\mathcal{D}, \hat{\mu}_S) &\leq \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta|S)} \frac{1}{n} \sum_{m=1}^n \mathbb{E}_X \mathbb{E}_{q(e_{(m)}|\mathbf{e}, \phi, S_m)} \|x - g_\theta(e_{(m)})\|^2 \\ &= \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta|S)} \frac{1}{n} \sum_{m=1}^n \sum_{k=1}^K \|x - g_\theta(e_k)\|^2 \mathbb{E}_{q(J_m|\mathbf{e}, \phi, S_m)} \mathbb{1}_{k=J_m}. \end{aligned} \quad (22)$$

606 The first inequality is obtained by the definition of the Wasserstein distance.

607 The expression inside the square root corresponds to the first term of Eq.(9). We can verify this by
 608 noting that Eq.(22) represents the squared error at the test data point x under the prediction e_k , which

609 is derived using the training dataset. Meanwhile, the first term of Eq. (9) represents this error when
 610 the test data is replaced by the supersample \bar{U} .

611 Therefore, Eq.(22) can be upper-bounded by applying Eq.(13), which serves as the upper bound for
 612 Eq. (9).

$$\begin{aligned} & \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} W_2^2(\mathcal{D}, \hat{\mu}_S) \\ & \leq \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(e_{(m)} | \mathbf{e}, \phi, S_m)} \|S_m - g_\theta(e_{(m)})\|^2 + \frac{1}{\lambda} \text{KL}(\mathbf{Q} | \mathbf{P}) + \frac{\lambda \Delta^2}{2n} + \frac{\Delta}{\sqrt{n}}, \end{aligned} \quad (23)$$

613 where

$$\begin{aligned} \mathbf{Q} & := q(\mathbf{e}, \phi, \theta | S) \prod_{m=1}^n q(J_m | \mathbf{e}, \phi, S_m) = q(\mathbf{e}, \phi, \theta | S) \prod_{m=1}^n q(e_{(m)} | \mathbf{e}, \phi, S_m), \\ \mathbf{P} & := q(\mathbf{e}, \phi, \theta | S) \prod_{m=1}^n q(J_m | \mathbf{e}, \phi) = q(\mathbf{e}, \phi, \theta | S) \prod_{m=1}^n q(e_{(m)} | \mathbf{e}, \phi). \end{aligned}$$

614 Next, the second term of Eq. (21) is bounded as follows; We express $\prod_{m=1}^n q(e_{(m)} | \mathbf{e}, \phi) = q(e)$ for
 615 simplicity, then we have

$$\begin{aligned} & \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} W_2^2(\hat{\mu}_S, \hat{\mu}) \\ & \leq \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(e)} \mathbb{E}_{q(e_{(m)} | \mathbf{e}, \phi, S_m)} \|g_\theta(e) - g_\theta(e_{(m)})\|^2 \\ & = \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(e)} \|g_\theta(e)\|^2 + \mathbb{E}_{q(e_{(m)} | \mathbf{e}, \phi, S_m)} \|g_\theta(e_{(m)})\|^2 - 2\mathbb{E}_{q(e)} g_\theta(e) \cdot \mathbb{E}_{q(e_{(m)} | \mathbf{e}, \phi, S_m)} g_\theta(e_{(m)}) \\ & \leq \frac{1}{\lambda} \text{KL}(\mathbf{Q} | \mathbf{P}) + \frac{\lambda \Delta^2}{2n}, \end{aligned} \quad (24)$$

616 where we used the Donsker Valadhan lemma for the first and third terms, changing the expectation
 617 from \mathbf{Q} to \mathbf{P} .

618 Combining Eqs. (23) and (24), we have

$$\begin{aligned} & \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} W_2^2(\mathcal{D}, \hat{\mu}) \\ & \leq 2 \left(\mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \frac{1}{n} \sum_{m=1}^n \mathbb{E}_{q(e_{(m)} | \mathbf{e}, \phi, S_m)} \|S_m - g_\theta(e_{(m)})\|^2 + \frac{2}{\lambda} \text{KL}(\mathbf{Q} | \mathbf{P}) + \frac{\lambda \Delta^2}{n} + \frac{\Delta}{\sqrt{n}} \right). \end{aligned}$$

619 Then by optimizing λ , we have

$$\mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} W_2^2(\mathcal{D}, \hat{\mu}) \leq \mathbb{E}_S \mathbb{E}_{q(\mathbf{e}, \phi, \theta | S)} \frac{2}{n} \sum_{m=1}^n \mathbb{E}_{q(e_{(m)} | \mathbf{e}, \phi, S_m)} \|S_m - g_\theta(e_{(m)})\|^2 + 2\Delta \sqrt{\frac{2}{n} \text{KL}(\mathbf{Q} | \mathbf{P})} + \frac{2\Delta}{\sqrt{n}}.$$